

## DELOCALIZATION FOR A CLASS OF RANDOM BLOCK BAND MATRICES

ZHIGANG BAO AND LÁSZLÓ ERDŐS

ABSTRACT. We consider  $N \times N$  Hermitian random matrices  $H$  consisting of blocks of size  $M \geq N^{6/7}$ . The matrix elements are i.i.d. within the blocks, close to a Gaussian in the four moment matching sense, but their distribution varies from block to block to form a block-band structure, with an essential band width  $M$ . We show that the entries of the Green's function  $G(z) = (H - z)^{-1}$  satisfy the local semicircle law with spectral parameter  $z = E + i\eta$  down to the real axis for any  $\eta \gg N^{-1}$ , using a combination of the supersymmetry method inspired by [19] and the Green's function comparison strategy. Previous estimates were valid only for  $\eta \gg M^{-1}$ . The new estimate also implies that the eigenvectors in the middle of the spectrum are fully delocalized.

## 1. INTRODUCTION

The Hamiltonian of quantum systems on a graph  $\Gamma$  is a self-adjoint matrix  $H = (h_{ab})_{a,b \in \Gamma}$ ,  $H = H^*$ . The matrix elements  $h_{ab}$  represent the quantum transition rates from vertex  $a$  to  $b$ . Disordered quantum systems have random matrix elements. We assume they are centered,  $\mathbb{E}h_{ab} = 0$ , and independent subject to the basic symmetry constraint  $h_{ab} = \bar{h}_{ba}$ . The variance  $\sigma_{ab}^2 := \mathbb{E}|h_{ab}|^2$  represents the strength of the transition from  $a$  to  $b$  and we use a scaling where the norm  $\|H\|$  is typically order 1. The simplest case is the mean field model, where  $h_{ab}$  are identically distributed; this is the standard Wigner matrix ensemble [29]. The other prominent example is the Anderson model [2] or random Schrödinger operator,  $H = \Delta + V$ , where the kinetic energy  $\Delta$  is the (deterministic) graph Laplacian and the potential  $V = (V_x)_{x \in \Gamma}$  is an on-site multiplication operator with random multipliers. If  $\Gamma$  is a discrete  $d$ -dimensional torus, then only few matrix elements  $h_{ab}$  are nonzero and they connect nearest neighbor points in the torus,  $\text{dist}(a, b) \leq 1$ . This is in sharp contrast to the mean field character of the Wigner matrices.

Random band matrices naturally interpolate between the mean field Wigner matrices and the short range Anderson model. They are characterized by a parameter  $M$ , called the *band width*, such that the matrix elements  $h_{ab}$  for  $\text{dist}(a, b) \geq M$  are zero or negligible. If  $M$  is comparable with the diameter  $N$  of the system then we are in the mean field regime, while  $M \sim 1$  corresponds to the short range model.

The Anderson model exhibits a metal-insulator phase transition: at high disorder the system is in the localized (insulator) regime, while at small disorder it is in the delocalized (metallic) regime, at least in  $d \geq 3$  dimensions and away from the spectral edges. The localized regime is characterized by exponentially decaying eigenfunctions and off diagonal decay of the Green's function, while in the complementary regime the eigenfunctions are supported in the whole physical space. In terms of the *localization length*  $\ell$ , the characteristic length scale of the decay, the localized regime corresponds to  $\ell \ll N$ , while in the delocalized regime  $\ell \sim N$ . Starting from the basic papers [1, 15], the localized regime is well understood, but the delocalized regime is still an open mathematical problem for the  $d$ -dimensional torus.

---

*Date:* March 26, 2015.

*2010 Mathematics Subject Classification.* 15B52, 82B44.

*Key words and phrases.* Random band matrix, supersymmetry, Green's function comparison, local semicircle law, delocalization.

Z.G. Bao was supported by ERC Advanced Grant RANMAT No.338804; L. Erdős was partially supported by ERC Advanced Grant RANMAT No.338804 .

Since the eigenvectors of the mean field Wigner matrices are always delocalized [13, 14], while the short range models are localized, by varying the parameter  $M$  in the random band matrix, one expects a (de)localization phase transition. Indeed, for  $d = 1$  it is conjectured (and supported by non rigorous supersymmetric calculations [16]) that the system is delocalized for broad bands,  $M \gg N^{1/2}$  and localized for  $M \ll N^{1/2}$ . The optimal power  $1/2$  has not yet been achieved from either sides. Localization has been shown for  $M \ll N^{1/8}$  in [22], while delocalization in a certain sense was proven for  $M \gg N^{4/5}$  in [11]. Interestingly, for a special Gaussian model even the sine kernel behavior of the 2-point correlation function of the characteristic polynomials could be proven down to the optimal band width  $M \gg N^{1/2}$ , see [18, 20]. Note that the sine kernel is consistent with the delocalization but does not imply it. We remark that our discussion concerns the bulk of the spectrum; the transition at the spectral edge is much better understood. In [24] it was shown that the edge spectrum follows the Tracy Widom distribution, characteristic to mean field model, for  $M \gg N^{5/6}$ , but it yields a different distribution for narrow bands,  $M \ll N^{5/6}$ .

Delocalization is closely related to estimates on the diagonal elements of the resolvent  $G(z) = (H - z)^{-1}$  at spectral parameters with small imaginary part  $\eta = \text{Im} z$ . Indeed, if  $G_{ii}(E + i\eta)$  is bounded for all  $i$  and all  $E \in \mathbb{R}$ , then each  $\ell^2$ -normalized eigenvector  $\mathbf{u}$  of  $H$  is delocalized on scale  $\eta^{-1}$  in a sense that  $\max_i |u_i|^2 \lesssim \eta$ , i.e.  $u$  is supported on at least  $\eta^{-1}$  sites. In particular, if  $G_{ii}$  can be controlled down to the scale  $\eta \sim 1/N$ , then the system is in the complete delocalized regime. Moreover, boundedness of  $G_{ii}$  also implies that the local semicircle law holds for the same regime of  $\eta$ .

For band matrices with band width  $M$ , or even under the more general condition  $\sigma_{ab}^2 \leq M^{-1}$ , the boundedness of  $G_{ii}$  was shown down to scale  $\eta \gg M^{-1}$  in [14] (see also [12]). If  $M \gg N^{1/2}$ , it is expected that  $G_{ii}$  remains bounded even down to  $\eta \gg N^{-1}$  which is the typical eigenvalue spacing, the smallest relevant scale in the model. However, with the standard approach [14, 12] via the self-consistent equations for the Green's function does not seem to work for  $\eta \leq 1/M$ ; the fluctuation is hard to control. The more subtle approach using the self-consistent *matrix* equation in [11] could prove delocalization and the off-diagonal Green's function profile that are consistent with the conventional quantum diffusion picture, but it was valid only for relatively large  $\eta$ , far from  $M^{-1}$ . Moment methods, even with a delicate renormalization scheme [23] could not break the barrier  $\eta \sim M^{-1}$  either.

In this paper we attack the problem differently; with supersymmetric (SUSY) techniques. Our main result is that  $G_{ii}(z)$  is bounded, and the local semicircle law holds for any  $\eta \gg N^{-1}$ , i.e. down to the optimal scale, if the band width is not too small,  $M \gg N^{6/7}$ , but under two technical assumptions. First, we consider a generalization of Wegner's  $n$ -orbital model [21, 28], namely, we assume that the band matrix has a block structure, i.e. it consists of  $M \times M$  blocks and the matrix elements within each block have the same distribution. This assumption is essential to reduce the number of integration variables in the supersymmetric representation, since, roughly speaking, each  $M \times M$  block will be represented by a single supermatrix with 16 supersymmetric variables. Second, we assume that the distribution of the matrix elements matches a Gaussian up to four moments in the spirit of [27]. Supersymmetry heavily uses Gaussian integrations, in fact all mathematically rigorous works on random band matrices with supersymmetric method assume that the matrix elements are Gaussian, see [4, 5, 6, 18, 19, 20, 25, 26]. The Green's function comparison method [14] allows one to compare Green's functions of two matrix ensembles provided that the distributions match up to four moments and provided that  $G_{ii}$  are bounded. This was an important motivation to reach the optimal scale  $\eta \gg N^{-1}$ .

In the next subsections we introduce the model precisely and state our main results. While SUSY approach is ubiquitous in physics, see e.g. the basic monograph by Efetov [7], its application in rigorous proofs is notoriously difficult. Initiated by T. Spencer (see [25] for a summary) and starting with the

paper [4] by Disertori, Pinsky and Spencer, only a handful of mathematical papers have succeeded in exploiting this powerful tool. Our supersymmetric analysis was inspired by [19], but our observable,  $G_{ab}$ , requires a partly different formalism, in particular we use the singular version of the superbosonization formula [3]. Moreover, our analysis is considerably more involved since we consider relatively narrow bands. In Section 1.3, we explain our novelties compared with [19].

**1.1. Matrix model.** Let

$$H_N = (h_{ab})$$

be an  $N \times N$  random Hermitian matrix, in which the entries are independent (up to symmetry), centered, complex variables. In this paper, we are concerned with  $H_N$  possessing a block band structure. To define this structure explicitly, we set the additional parameters  $M \equiv M(N)$  and  $W \equiv W(N)$  satisfying

$$W = N/M.$$

For simplicity, we assume that both  $M$  and  $W$  are integers. Let  $S = (\mathfrak{s}_{jk})$  be a  $W \times W$  symmetric matrix, which will be chosen as a weighted Laplacian of a connected graph on  $W$  vertices. Now, we decompose  $H_N$  into  $W \times W$  blocks of size  $M \times M$ , and relabel

$$h_{jk,\alpha\beta} := h_{ab}, \quad j, k = 1, \dots, W, \quad \alpha, \beta = 1, \dots, M,$$

where  $(j, k) = (\lceil a/M \rceil, \lceil b/M \rceil)$  is the index of the block containing  $h_{ab}$ , and

$$(\alpha, \beta) = (a - (j - 1)M, b - (k - 1)M)$$

describes the location of the entry in the block. Moreover, we assume

$$\mathbb{E} h_{jk,\alpha\beta} h_{j'k',\alpha'\beta'} = \frac{1}{M} \delta_{jk'} \delta_{j'k} \delta_{\alpha\beta'} \delta_{\beta\alpha'} (\delta_{jk} + \mathfrak{s}_{jk}). \quad (1.1)$$

That means, the variance profile of the random matrix  $\sqrt{M}H_N$  is given by

$$\tilde{S} = (\tilde{\mathfrak{s}}_{jk}) := I + S, \quad (1.2)$$

in which each entry represents the common variance of the entries in the corresponding block of  $\sqrt{M}H_N$ . Moreover, if  $h_{jk,\alpha\beta}$ 's are Gaussian, (1.1) also implies that for each off-diagonal entry  $h_{jk,\alpha\beta}$ , its real part and imaginary part are i.i.d.  $N(0, \tilde{\mathfrak{s}}_{jk}/2M)$  variables.

**1.2. Assumptions and main results.** In the sequel, for some matrix  $A = (a_{ij})$  and some index sets  $\mathbf{I}$  and  $\mathbf{J}$ , we introduce the notation  $A^{(\mathbf{I}|\mathbf{J})}$  to denote the submatrix obtained by deleting the  $i$ -th row and  $j$ -th column of  $A$  for all  $i \in \mathbf{I}$  and  $j \in \mathbf{J}$ . We will adopt the abbreviation

$$A^{(i|j)} := A^{(\{i\}|\{j\})}, \quad i \neq j, \quad A^{(i)} := A^{(\{i\}|\{i\})}. \quad (1.3)$$

In addition, we use  $\|A\|_{\max} := \max_{i,j} |a_{ij}|$  to denote the max norm of  $A$ . Throughout the paper, we need some assumptions on  $S$ .

**Assumption 1.1** (On  $S$ ). *Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a connected simple graph with  $\mathcal{V} = \{1, \dots, W\}$ . Assume that  $S$  is a  $W \times W$  symmetric matrix satisfying the following four conditions.*

- (i)  *$S$  is a weighted Laplacian on  $\mathcal{G}$ , i.e. for  $i \neq j$ , we have  $\mathfrak{s}_{ij} > 0$  if  $\{i, j\} \in \mathcal{E}$  and  $\mathfrak{s}_{ij} = 0$  if  $\{i, j\} \notin \mathcal{E}$ , and for the diagonal entries, we have*

$$\mathfrak{s}_{ii} = - \sum_{j:j \neq i} \mathfrak{s}_{ij}, \quad \forall i = 1, \dots, W.$$

- (ii)  $\tilde{S}$  defined in (1.2) is strictly diagonally dominant, i.e., there exists some constant  $c_0 > 0$  such that

$$1 + 2\mathfrak{s}_{ii} > c_0, \quad \forall i = 1, \dots, W.$$

- (iii) For the discrete Green's functions, we assume that there exist some positive constants  $C$  and  $\gamma$  such that

$$\max_{i=1, \dots, W} \|(S^{(i)})^{-1}\|_{\max} \leq CW^\gamma.$$

- (iv) There exists a spanning tree  $\mathcal{G}_0 = (\mathcal{V}, \mathcal{E}_0) \subset \mathcal{G}$ , on which the weights are bounded below, i.e. for some constant  $c > 0$ , we have

$$\mathfrak{s}_{ij} \geq c, \quad \text{if } \{i, j\} \in \mathcal{E}_0$$

**Remark 1.2.** From Assumption 1.1 (ii), we easily see that

$$\tilde{S} \geq c_0 I \tag{1.4}$$

for the same positive constant  $c_0$ . In addition, the lower bound  $c$  in (iv) can be weakened to  $N^{-\varepsilon}$  for some sufficiently small constant  $\varepsilon > 0$ . But for simplicity, we will not try to optimize this bound in this paper.

**Example 1.1.** Let  $\Delta$  be the standard discrete Laplacian on the  $\mathbf{d}$ -dimensional torus  $[1, \mathfrak{w}]^{\mathbf{d}} \cap \mathbb{Z}^{\mathbf{d}}$ , with periodic boundary condition, where  $\mathfrak{w} = W/\mathbf{d}$ . Here by standard we mean the weights on the edges of the box are all 1. Now let  $S = a\Delta$  for some positive constant  $a < 1/4\mathbf{d}$ . It is then easy to check Assumption 1.1 (i), (ii) and (iv) are satisfied. In addition, if  $\mathbf{d} = 1$ , it is well known that we can choose  $\gamma = 1$  in Assumption 1.1 (iii). For  $\mathbf{d} \geq 3$ , one can choose  $\gamma = 0$ . For  $\mathbf{d} = 2$ , one can choose  $\gamma = \varepsilon$  for arbitrarily small constant  $\varepsilon$ . For instance, one can refer to [8] for more details.

For simplicity, we also introduce the notation

$$\sigma_{ij}^2 := \mathbb{E}|h_{ij}|^2, \quad \mathcal{T} := (\sigma_{ij}^2)_{N,N} = \frac{1}{M} \tilde{S} \otimes \mathbf{1}_M \mathbf{1}_M', \quad i, j = 1, \dots, N, \tag{1.5}$$

where  $\mathbf{1}_M$  is the  $M$ -dimensional vector whose components are all 1 and  $\tilde{S}$  is the variance matrix in (1.2). It is elementary that

$$\text{Spec}(\mathcal{T}) = \text{Spec}(\tilde{S}) \cup \{0\} \subset [0, 1]. \tag{1.6}$$

Our assumption on  $M$  depends on the constant  $\gamma$  in Assumption 1.1 (iii).

**Assumption 1.3** (On  $M$ ). We assume that there exists a (small) positive constant  $\varepsilon_1$  such that

$$M \geq W^{4+2\gamma+\varepsilon_1}. \tag{1.7}$$

**Remark 1.4.** A direct consequence of (1.7) and  $N = MW$  is

$$M \geq N^{\frac{4+2\gamma+\varepsilon_1}{5+2\gamma+\varepsilon_1}}. \tag{1.8}$$

Especially, when  $\gamma = 1$ , one has  $M \gg N^{6/7}$ . Actually, through a more involved analysis, (1.7) (or (1.8)) can be further improved. At least, for  $\gamma \leq 1$ , we expect that  $M \gg N^{4/5}$  is enough. However, we will not pursue this direction here.

Besides Assumption 1.1 on the variance profile of  $H$ , we need to impose some additional assumption on the distribution of its entries. To this end, we temporarily employ the notation  $H^g = (h_{ab}^g)$  to represent a random block band matrix with Gaussian entries, satisfying (1.1), Assumption 1.1 and Assumption 1.3.

**Assumption 1.5** (On distribution). *We assume that for each  $a, b \in \{1, \dots, N\}$ , the moments of the entry  $h_{ab}$  match those of  $h_{ab}^g$  up to the 4th order, i.e.*

$$\mathbb{E}(\operatorname{Re} h_{ab})^k (\operatorname{Im} h_{ab})^\ell = \mathbb{E}(\operatorname{Re} h_{ab}^g)^k (\operatorname{Im} h_{ab}^g)^\ell, \quad \forall k, \ell \in \mathbb{N}, \quad \text{s.t. } k + \ell \leq 4. \quad (1.9)$$

*In addition, we assume the distribution of  $h_{ab}$  possesses a subexponential tail, namely, there exist positive constants  $c_1$  and  $c_2$  such that for any  $\tilde{\gamma} > 0$ ,*

$$\mathbb{P}\left(|h_{ab}| \geq \tilde{\gamma}^{c_1} (\mathbb{E}|h_{ab}|^2)^{\frac{1}{2}}\right) \leq c_2 e^{-\tilde{\gamma}} \quad (1.10)$$

*holds uniformly for all  $a, b = 1, \dots, N$ .*

The four moment condition (1.9) in the context of random matrices first appeared in Tao and Vu's work [27].

To state our results, we will need the following notion on the comparison of two random sequences, which was introduced in [9] and [12].

**Definition 1.6** (Stochastic domination). *For some possibly  $N$ -dependent parameter set  $\mathbf{U}_N$ , and two families of random variables  $\mathbf{X} = (\mathbf{X}_N(u) : N \in \mathbb{N}, u \in \mathbf{U}_N)$  and  $\mathbf{Y} = (\mathbf{Y}_N(u) : N \in \mathbb{N}, u \in \mathbf{U}_N)$ , we say that  $\mathbf{X}$  is stochastically dominated by  $\mathbf{Y}$ , if for all  $\varepsilon > 0$  and  $D > 0$  we have*

$$\sup_{u \in \mathbf{U}_N} \mathbb{P}\left(\mathbf{X}_N(u) \geq N^\varepsilon \mathbf{Y}_N(u)\right) \leq N^{-D} \quad (1.11)$$

*for all sufficiently large  $N \geq N_0(\varepsilon, D)$ . In this case we write*

$$\mathbf{X} \prec \mathbf{Y}.$$

For example, by (1.1) and Assumption 1.5, we have

$$|h_{ab}| \prec \frac{1}{\sqrt{M}}, \quad \forall a, b = 1, \dots, N. \quad (1.12)$$

Note that  $\tilde{S}$  is doubly stochastic. It is known that the empirical eigenvalue distribution of  $H_N$  converges to the semicircle law, whose density function is given by

$$\varrho_{sc}(x) := \frac{1}{2\pi} \sqrt{4 - x^2} \cdot \mathbf{1}(|x| \leq 2).$$

We denote the Green's function of  $H_N$  by

$$G(z) \equiv G_N(z) := (H_N - z)^{-1}, \quad z = E + i\eta \in \mathbb{C}^+ := \{w \in \mathbb{C} : \operatorname{Im} w > 0\}$$

and its  $(a, b)$  matrix element is  $G_{ab}(z)$ . Throughout the paper, we will always use  $E$  and  $\eta$  to denote the real and imaginary part of  $z$  without further mention. In addition, for simplicity, we suppress the subscript  $N$  from the notation of the matrices here and there. The Stieltjes transform of  $\varrho_{sc}(x)$  is

$$m_{sc}(z) = \int_{-2}^2 \frac{\varrho_{sc}(x)}{x - z} dx = \frac{-z + \sqrt{z^2 - 4}}{2},$$

where we chose the branch of the square root with positive imaginary part for  $z \in \mathbb{C}^+$ . Note that  $m_{sc}(z)$  is a solution to the following self-consistent equation

$$m_{sc}(z) = \frac{1}{-z - m_{sc}(z)}. \quad (1.13)$$

The semicircle law also holds in a local sense, see Theorem 2.3 in [12]. For simplicity, we cite this result with a slight modification adjusted to our assumption.

**Proposition 1.7** (Erdős, Knowles, Yau, Yin, [12]). *Let  $H$  be a random block band matrix satisfying Assumptions 1.1, 1.3 and 1.5. Then*

$$\max_{a,b} |G_{ab}(z) - \delta_{ab} m_{sc}(z)| \prec \frac{1}{\sqrt{M\eta}}, \quad \text{if } E \in [-2 + \kappa, 2 - \kappa] \quad \text{and} \quad M^{-1+\varepsilon} \leq \eta \leq 10 \quad (1.14)$$

for any fixed small positive constants  $\kappa$  and  $\varepsilon$ .

**Remark 1.8.** We remark that Theorem 2.3 in [12] was established under a more general assumption  $\sum_k \sigma_{jk}^2 = 1$  and  $\sigma_{jk}^2 \leq C/M$ . Especially, the block structure on the variance profile is not needed. In addition, Theorem 2.3 in [12] also covers the edges of the spectrum, which will not be discussed in this paper. We also refer to [14] for a previous result, see Theorem 2.1 therein.

Our aim in this paper is to extend the local semicircle law to the regime  $\eta \gg N^{-1}$  and replace  $M$  with  $N$  in (1.14). More specifically, we will work in the following set, defined for arbitrarily small constant  $\kappa > 0$  and any sufficiently small positive constant  $\varepsilon_2 := \varepsilon_2(\varepsilon_1)$ ,

$$\mathbf{D}(N, \kappa, \varepsilon_2) := \left\{ z = E + i\eta \in \mathbb{C} : |E| \leq \sqrt{2} - \kappa, N^{-1+\varepsilon_2} \leq \eta \leq M^{-1} N^{\varepsilon_2} \right\}. \quad (1.15)$$

Throughout the paper, we will assume that  $\varepsilon_2$  is much smaller than  $\varepsilon_1$ , see (1.7) for the latter. Specifically, there exists some large enough constant  $C$  such that  $\varepsilon_2 \leq \varepsilon_1/C$ .

**Theorem 1.9** (Local semicircle law). *Suppose that  $H$  is a random block band matrix satisfying Assumptions 1.1, 1.3 and 1.5. Let  $\kappa$  be an arbitrarily small positive constant and  $\varepsilon_2$  be any sufficiently small positive constant. Then*

$$\max_{a,b} |G_{ab}(z) - \delta_{ab} m_{sc}(z)| \prec \frac{1}{\sqrt{N\eta}} \quad (1.16)$$

holds uniformly on  $\mathbf{D}(N, \kappa, \varepsilon_2)$ .

**Remark 1.10.** The restriction  $|E| \leq \sqrt{2} - \kappa$  in (1.15) is technical. We believe the result can be extended to the whole bulk regime of the spectrum, i.e.,  $|E| \leq 2 - \kappa$ , see Section 12 for further comment. The upper bound of  $\eta$  in (1.15) is also technical. However, for  $\eta > M^{-1} N^{\varepsilon_2}$ , one can control the Green's function by (1.14) directly.

Let  $\lambda_1, \dots, \lambda_N$  be the eigenvalues of  $H_N$ . We denote by  $\mathbf{u}_i := (u_{i1}, \dots, u_{iN})$  the normalized eigenvector of  $H_N$  corresponding to  $\lambda_i$ . From Theorem 1.9, we can also get the following delocalization property for the eigenvectors.

**Theorem 1.11** (Complete delocalization). *Let  $H$  be a random block band matrix satisfying Assumptions 1.1, 1.3 and 1.5. We have*

$$\max_{i: |\lambda_i| \leq \sqrt{2} - \kappa} \|\mathbf{u}_i\|_\infty \prec N^{-\frac{1}{2}}. \quad (1.17)$$

**Remark 1.12.** We remark that delocalization in a certain weak sense was proven in [11] for an even more general class of random band matrices if  $M \gg N^{4/5}$ . However, Theorem 1.11 asserts delocalization for **all** eigenvectors in a very strong sense (supremum norm), while Proposition 7.1 of [11] stated that **most** eigenvectors are delocalized in a sense that their substantial support cannot be too small.

**1.3. Outline of the proof strategy and novelties.** In this section, we briefly outline the strategy for the proof of Theorem 1.9.

The first step, which is the main task of the whole proof, is to establish the following Theorem 1.14, namely, a prior estimate of the Green's function in the Gaussian case. For technical reason, we need the following slight modification of Assumption 1.3, to state the result.

**Assumption 1.13** (On  $M$ ). Let  $\varepsilon_1$  be the small positive constant in Assumption 1.3. We assume

$$N(\log N)^{-10} \geq M \geq W^{4+2\gamma+\varepsilon_1}. \quad (1.18)$$

In the regime  $M \geq N(\log N)^{-10}$ , we see that (1.16) anyway follows from (1.14) directly.

**Theorem 1.14.** Assume that  $H$  is a Gaussian block band matrix, satisfying Assumptions 1.1 and 1.13. Let  $n$  be any fixed positive integer. Let  $\kappa$  be an arbitrarily small positive constant and  $\varepsilon_2$  be any sufficiently small positive constant. There is  $N_0 = N_0(n)$ , such that for all  $N \geq N_0$  and all  $z \in \mathbf{D}(N, \kappa, \varepsilon_2)$ , we have

$$\mathbb{E}|G_{ab}(z)|^{2n} \leq N^{C_0} \left( \delta_{ab} + \frac{1}{(N\eta)^n} \right), \quad \forall a, b = 1, \dots, N \quad (1.19)$$

for some positive constant  $C_0$  independent of  $n$  and  $z$ .

**Remark 1.15.** Much more delicate analysis can show that the prefactor  $N^{C_0}$  can be improved to some  $n$ -dependent constant  $C_n$ . We refer to Section 12 for further comment on this issue.

Using the definition of stochastic domination in Definition 1.6, a simple Markov inequality shows that (1.19) implies

$$|G_{ab}(z)| \prec \delta_{ab} + \frac{1}{\sqrt{N\eta}}, \quad \forall a, b = 1, \dots, N. \quad (1.20)$$

The proof of Theorem 1.14 is the main task of our paper. We will use the supersymmetry method. We partially rely on the arguments from Shcherbina's work [19] concerning universality of the local 2-point function and we develop new techniques to treat our observable, the high moment of the entries of  $G(z)$ , under a more general setting. We will comment on the novelties later in this subsection.

The second step is to generalize Theorem 1.14 from the Gaussian case to more general distribution satisfying Assumption 1.5, via a Green's function comparison strategy initiated in [14], see Lemma 2.1 below.

The last step is to use Lemma 2.1 and its Corollary 2.2 to prove our main theorems. Using (1.20) below to bound the error term in the self-consistent equation for the Green's function, we can prove Theorem 1.9 by a continuity argument in  $z$ , with the aid of the initial estimate for large  $\eta$  provided in Proposition 1.7. Theorem 1.11 will then easily follow from Theorem 1.9.

The second and the last steps are carried out in Section 2. The main body of this paper, Sections 3–11 is devoted to the proof of Theorem 1.14.

One of the main novelty of this work is to combine the supersymmetry method and the Green's function comparison strategy to go beyond the Gaussian ensemble, which was so far the only random band matrix ensemble amenable to the supersymmetry method, as mentioned at the beginning. The comparison strategy requires an apriori control on the individual matrix elements of the Green's function with high probability (see (1.20)), this is one of our main motivations behind Theorem 1.14.

Although we consider a different observable than [19], many technical aspects of the supersymmetric analysis overlaps with [19]. For the convenience of the reader, we now briefly introduce the strategy of [19], and highlight the main novelties of our work.

In [19], the author considers the 2-point correlation function of the trace of the resolvent of the Gaussian block band matrix  $H$ , with the variance profile  $\tilde{S} = 1 + a\Delta$ , under the assumption  $M \sim N$  (note that we use  $M$  instead of  $W$  in [19] for the size of the blocks). The 2-point correlation function can be expressed in terms of a superintegral of a superfunction  $F(\{\check{\mathcal{S}}_i\}_{i=1}^W)$  with a collection of  $4 \times 4$  supermatrices  $\check{\mathcal{S}}_i := \mathcal{Z}_i^* \mathcal{Z}_i$ . Here for each  $i$ ,  $\mathcal{Z}_i = (\Psi_{1,i}, \Psi_{2,i}, \Phi_{1,i}, \Phi_{2,i})$  is an  $M \times 4$  matrix and  $\mathcal{Z}_i^*$  is its conjugate transpose, where  $\Psi_{1,i}$  and  $\Psi_{2,i}$  are Grassmann  $M$ -vectors whilst  $\Phi_{1,i}$  and  $\Phi_{2,i}$  are complex  $M$ -vectors. Then, by using the superbosonization formula in the nonsingular case ( $M \geq 4$ ) from [17],



one can transform the superintegral of  $F(\{\check{\mathcal{S}}_i\}_{i=1}^W)$  to a superintegral of  $F(\{\mathcal{S}_i\}_{i=1}^W)$ , where each  $\mathcal{S}_i$  is a supermatrix akin to  $\check{\mathcal{S}}_i$ , but only consists of 16 independent variables (either complex or Grassmann). We will call the integral representation of the observable after using the superbosonization formula as the *final integral representation*. Schematically it has the form

$$\int \mathbf{g}(\mathcal{S}_c) e^{M\mathbf{f}_c(\mathcal{S}_c) + \mathbf{f}_g(\mathcal{S}_g, \mathcal{S}_c)} d\mathcal{S}, \quad (1.21)$$

for some functions  $\mathbf{g}(\cdot)$ ,  $\mathbf{f}_c(\cdot)$  and  $\mathbf{f}_g(\cdot)$ , where we used the abbreviation  $\mathcal{S} := \{\mathcal{S}_i\}_{i=1}^W$  and  $\mathcal{S}_c$  and  $\mathcal{S}_g$  represents the collection of all complex variables and Grassmann variables in  $\mathcal{S}$ , respectively. Here,  $\mathbf{g}(\mathcal{S}_c)$  and  $\mathbf{f}_c(\mathcal{S}_c)$  are some complex functions and  $\mathbf{f}_g(\mathcal{S}_g, \mathcal{S}_c)$  will be mostly regarded as a function of the Grassmann variables with complex variables as its parameters. The number of variables (either complex or Grassmann) in the final integral representation then turns out to be of order  $W$ , which is much smaller than the original order  $N$ . In fact, in [19] it is assumed that  $W = O(1)$  although the author also mentions the possibility to deal with the case  $W \sim N^\varepsilon$  for some small positive  $\varepsilon$ , see the remark below Theorem 1 therein.

Performing a saddle point analysis for the complex measure  $\exp\{M\mathbf{f}_c(\mathcal{S}_c)\}$ , one can restrict the integral in a small vicinity of some saddle point, say,  $\mathcal{S}_c = \mathcal{S}_{c0}$ . It turns out that  $\mathbf{f}_c(\mathcal{S}_{c0}) = 0$  and  $\mathbf{f}_c(\mathcal{S}_c)$  decays quadratically away from  $\mathcal{S}_{c0}$ . Consequently, by plugging in the saddle point  $\mathcal{S}_{c0}$ , one can estimate  $\mathbf{g}(\mathcal{S}_c)$  by  $\mathbf{g}(\mathcal{S}_{c0})$  directly. However, for  $\exp\{M\mathbf{f}_c(\mathcal{S}_c)\}$  and  $\exp\{\mathbf{f}_g(\mathcal{S}_g, \mathcal{S}_c)\}$ , one shall expand them around the saddle point. Roughly speaking, in some vicinity of  $\mathcal{S}_{c0}$ , one will find that the expansions read

$$e^{M\mathbf{f}_c(\mathcal{S}_c)} = \exp\{-\mathbf{u}'\mathbb{A}\mathbf{u} + \mathbf{e}_c(\mathbf{u})\}, \quad e^{\mathbf{f}_g(\mathcal{S}_g, \mathcal{S}_c)} = \exp\{-\boldsymbol{\rho}'\mathbb{H}\boldsymbol{\tau}\} \mathbf{p}(\boldsymbol{\rho}, \boldsymbol{\tau}, \mathbf{u}), \quad (1.22)$$

where  $\mathbf{u}$  is a complex vector of dimension  $O(W)$ , which is essentially a vectorization of  $\sqrt{M}(\mathcal{S}_c - \mathcal{S}_{c0})$ ;  $\mathbf{e}_c(\mathbf{u}) = o(1)$  is some error term;  $\boldsymbol{\rho}$  and  $\boldsymbol{\tau}$  are two Grassmann vectors of dimension  $O(W)$ ;  $\mathbb{A}$  is a complex matrix with positive-definite Hermitian part and  $\mathbb{H}$  is a complex matrix;  $\mathbf{p}(\boldsymbol{\rho}, \boldsymbol{\tau}, \mathbf{u})$  is the expansion of  $\exp\{\mathbf{f}_g(\mathcal{S}_g, \mathcal{S}_c) - \mathbf{f}_g(\mathcal{S}_g, \mathcal{S}_{c0})\}$ , which possesses the form

$$\mathbf{p}(\boldsymbol{\rho}, \boldsymbol{\tau}, \mathbf{u}) = \sum_{\ell=0}^{O(W)} M^{-\frac{\ell}{2}} \mathbf{p}_\ell(\boldsymbol{\rho}, \boldsymbol{\tau}, \mathbf{u}), \quad (1.23)$$

where  $\mathbf{p}_\ell(\boldsymbol{\rho}, \boldsymbol{\tau}, \mathbf{u})$  is a polynomial of the components of  $\boldsymbol{\rho}$  and  $\boldsymbol{\tau}$  with degree  $2\ell$ , regarding  $\mathbf{u}$  as fixed parameters. Now, keeping the leading order term of  $\mathbf{p}(\boldsymbol{\rho}, \boldsymbol{\tau}, \mathbf{u})$ , and discarding the remainder terms, one can get the final estimate of the integral by taking the Gaussian integral over  $\mathbf{u}$ ,  $\boldsymbol{\rho}$  and  $\boldsymbol{\tau}$ . This completes the summary of [19].

Similarly to [19], we also use the superbosonization formula to reduce the number of variables and perform the saddle point analysis on the resulting integral. However, owing to the following three main aspects, our analysis is significantly different from [19].

•(Different observable) Our objective is to compute high moments of the single entry of the Green's function. By using Wick's formula (see Proposition 3.1), we express  $\mathbb{E}|G_{jk}|^{2n}$  in terms of a superintegral of some superfunction of the form

$$\tilde{F}\left(\{\Psi_{a,j}, \Psi_{a,j}^*, \Phi_{a,j}, \Phi_{a,j}^*\}_{a=1,2; j=1,\dots,W}\right) := (\bar{\phi}_{1,q,\beta} \phi_{1,p,\alpha} \bar{\phi}_{2,p,\alpha} \phi_{2,q,\beta})^n F(\{\check{\mathcal{S}}_i\}_{i=1}^W)$$

for some  $p, q \in \{1, \dots, W\}$  and  $\alpha, \beta \in \{1, \dots, M\}$ , where  $\phi_{1,p,\alpha}$  is the  $\alpha$ -th coordinate of  $\Phi_{1,p}$ , and the others are defined analogously. Unlike the case in [19],  $\tilde{F}$  is not a function of  $\{\check{\mathcal{S}}_i\}_{i=1}^W$  only. Hence, using the superbosonization formula to change  $\check{\mathcal{S}}_i$  to  $\mathcal{S}_i$  directly is not feasible in our case. In order to handle the factor  $(\bar{\phi}_{1,q,\beta} \phi_{1,p,\alpha} \bar{\phi}_{2,p,\alpha} \phi_{2,q,\beta})^n$ , the main idea is to split off certain rank-one supermatrices from  $\check{\mathcal{S}}_p$  and  $\check{\mathcal{S}}_q$  such that this factor can be expressed in terms of the entries of these rank-one supermatrices.



Then we use the superbosonization formula not only in the nonsingular case from [17] but also in the singular case from [3] to change and reduce the variables, resulting the final integral representation of  $\mathbb{E}|G_{jk}|^{2n}$ . Though this final integral representation, very schematically, is still of the form (1.21), due to the decomposition of the supermatrices  $\check{\mathcal{S}}_p$  and  $\check{\mathcal{S}}_q$ , it is considerably more complicated than its counterpart in [19]. Especially, the function  $\mathbf{g}(\mathcal{S}_c)$  differs from its counterpart in [19], and its estimate at the saddle point follows from a different argument.

•(Small band width) In [19], the author considers the case that the band width  $M$  is comparable with  $N$ , i.e. the number of blocks  $W$  is finite. Though the derivation of the 2-point correlation function is highly nontrivial even with such a large band width, our objective, the local semicircle law and delocalization of the eigenvectors, however, can be proved for the case  $M \sim N$  in a similar manner as for the Wigner matrix ( $M = N$ ), see [12, 14]. In our work, we will work with much smaller band width to go beyond the results in [12, 14], see Assumption 1.3. Several main difficulties stemming from a narrow band width can be heuristically explained as follows.

At first, let us focus on the integral over the small vicinity of the saddle point, in which the exponential functions in the integrand in (1.21) approximately look like (1.22).

We regard the first term in (1.22) as a complex Gaussian measure, of dimension  $O(W)$ . When  $W \sim 1$ , one can discard the error term  $\mathbf{e}_c(\mathbf{u})$  directly and perform the Gaussian integral over  $\mathbf{u}$ , due to the fact  $\int d\mathbf{u} \exp\{-\mathbf{u}'\text{Re}(\mathbb{A})\mathbf{u}\}|\mathbf{e}_c(\mathbf{u})| = o(1)$ . However, such an estimate is not allowed when  $W \sim N^\varepsilon$  (say), because the normalization of the measure  $\exp\{-\mathbf{u}'\text{Re}(\mathbb{A})\mathbf{u}\}$  might be exponentially larger than that of  $\exp\{-\mathbf{u}'\mathbb{A}\mathbf{u}\}$ . In order to handle this issue, we shall do a second deformation of the contours of the complex variables in the vicinity of the saddle, following the steepest descent paths exactly, whereby we can transform the complex Gaussian measure to a real one, thus the error term of the integral can be controlled.

Now, we turn to the second term in (1.22). When  $W \sim 1$ , there are only finitely many Grassmann variables. Hence, the complex coefficient of each term in the polynomial  $\mathbf{p}(\boldsymbol{\rho}, \boldsymbol{\tau}, \mathbf{u})$ , which is of order  $M^{-\ell/2}$  for some  $\ell \in \mathbb{N}$  (see (1.23)), actually controls the magnitude of the integral of this term against the Gaussian measure  $\exp\{-\boldsymbol{\rho}'\mathbb{H}\boldsymbol{\tau}\}$ . Consequently, in case of  $W \sim 1$ , it suffices to keep the leading order term (according to  $M^{-\ell/2}$ ), one may discard the others trivially, and compute the Gaussian integral over  $\boldsymbol{\rho}$  and  $\boldsymbol{\tau}$  explicitly. However, when  $W \sim N^\varepsilon$  (say), in light of the Wick's formula (3.2) and the fact that the coefficients are of order  $M^{-\ell/2}$ , the order of the integral of each term of  $\mathbf{p}(\boldsymbol{\rho}, \boldsymbol{\tau}, \mathbf{u})$  against the Gaussian measure reads  $M^{-\ell/2} \det \mathbb{H}^{(\mathbf{l}|\mathbf{j})}$  for some index sets  $\mathbf{l}$  and  $\mathbf{j}$  and some  $\ell \in \mathbb{N}$ . Due to the fact  $W \sim N^\varepsilon$ ,  $\det \mathbb{H}^{(\mathbf{l}|\mathbf{j})}$  is typically exponential in  $W$ . Hence, it is much more complicated to determine and compare the orders of the integrals of all  $e^{O(W)}$  terms. In our discussion, we perform a unified estimate for the integrals of all the terms, rather than simply compare them by  $M^{-\ell/2}$ .

In addition, the analysis for the integral away from the vicinity of the saddle point in our work is also quite different from [19]. Actually, the integral over the complement of the vicinity can be trivially ignored in [19], since each factor in the integrand of (1.21) is of order 1, thus gaining any  $o(1)$  factor for the integrand outside the vicinity is enough for the estimate. However, in our case, either  $\exp\{M\mathbf{f}_c(\mathcal{S}_c)\}$  or  $\int d\mathcal{S}_g \exp\{\mathbf{f}_g(\mathcal{S}_g, \mathcal{S}_c)\}$  is essentially exponential in  $W$ . This fact forces us to provide an a priori bound for  $\int d\mathcal{S}_g \exp\{\mathbf{f}_g(\mathcal{S}_g, \mathcal{S}_c)\}$  in the full domain of  $\mathcal{S}_c$  rather than in the vicinity of the saddle point only. In addition, an analysis of the tail behavior of the measure  $\exp\{M\mathbf{f}_c(\mathcal{S}_c)\}$  needs also to be performed.

•(General variance profile  $\tilde{S}$ ) In [19], the authors considered the special case  $S = a\Delta$  with  $a < 1/4d$ . We generalize the discussion to more general weighted Laplacians  $S$  satisfying Assumption 1.1, which, as a special case, includes the standard Laplacian  $\Delta$  for any fixed dimension  $d$ .

**1.4. Notation and organization.** Throughout the paper, we will need some notation. At first, we conventionally use  $U(r)$  to denote the unitary group of degree  $r$ , as well,  $U(1,1)$  represents the  $U(1,1)$  group. Furthermore, we denote

$$\mathring{U}(r) = U(r)/U(1)^r, \quad \mathring{U}(1,1) = U(1,1)/U(1)^2.$$

Recalling the real part  $E$  of  $z$ , we will frequently need the following two parameters

$$a_+ = \frac{\mathbf{i}E + \sqrt{4 - E^2}}{2}, \quad a_- = \frac{\mathbf{i}E - \sqrt{4 - E^2}}{2}.$$

Correspondingly, we define the following four matrices

$$D_{\pm} = \text{diag}(a_+, a_-), \quad D_{\mp} = \text{diag}(a_-, a_+), \quad D_+ = \text{diag}(a_+, a_+), \quad D_- = \text{diag}(a_-, a_-). \quad (1.24)$$

We remark here  $D_{\pm}$  does not mean “ $D_+$  or  $D_-$ ”. For simplicity, we introduce the following notation for some domains used throughout the paper.

$$\mathbb{I} := [0, 1], \quad \mathbb{L} := [0, 2\pi), \quad \Sigma : \text{unit circle}, \quad \mathbb{R}_+ := [0, \infty), \quad \mathbb{R}_- := -\mathbb{R}_+, \quad \Gamma := a_+\mathbb{R}_+.$$

For some  $\ell \times \ell$  Hermitian matrix  $A$ , we use  $\lambda_1(A) \leq \dots \leq \lambda_{\ell}(A)$  to represent its ordered eigenvalues. For some possibly  $N$ -dependent parameter set  $\mathbf{U}_N$ , and two families of complex functions  $\{a_N(u) : N \in \mathbb{N}, u \in \mathbf{U}_N\}$  and  $\{b_N(u) : N \in \mathbb{N}, u \in \mathbf{U}_N\}$ , if there exists a positive constant  $C > 1$  such that  $C^{-1}|b_N(u)| \leq |a_N(u)| \leq C|b_N(u)|$  holds uniformly in  $N$  and  $u$ , we write

$$a_N(u) \sim b_N(u).$$

Conventionally, we use  $\{\mathbf{e}_i : i = 1, \dots, \ell\}$  to denote the standard basis of  $\mathbb{R}^{\ell}$ , in which the dimension  $\ell$  has been suppressed for simplicity. For some real quantities  $a$  and  $b$ , we use  $a \wedge b$  and  $a \vee b$  to represent  $\min\{a, b\}$  and  $\max\{a, b\}$ , respectively.

Throughout the paper,  $c, c', c_1, c_2, C, C', C_1, C_2$  represent some generic positive constants that are possibly  $n$ -dependent and may differ from line to line. In contrast, we use  $C_0$  to denote some generic positive constant independent of  $n$ .

The paper will be organized in the following way. In Section 2, we prove Theorem 1.9 and Theorem 1.11, with Theorem 1.14. The proof of Theorem 1.14 will be done in Section 3–Section 11. More specifically, in Section 3, we use the supersymmetric formalism to represent  $\mathbb{E}|G_{ij}|^{2n}$  in terms of a superintegral, in which the integrand can be factorized into several functions; Section 4 is devoted to a preliminary analysis on these functions; Section 5–Section 10 are responsible for different steps of the saddle point analysis, whose organization will be further clarified at the end of Section 5; Section 11 is devoted to the final proof of Theorem 1.14, by summing up the discussions in 3–Section 10. Finally, in Section 12, we make some further comments on possible improvements.

## 2. PROOFS OF THEOREM 1.9 AND THEOREM 1.11

At first, (1.19) can be generalized to the generally distributed matrix with the four moment matching condition via the Green’s function comparison strategy.

**Lemma 2.1.** Assume that  $H$  is a random block band matrix, satisfying Assumptions 1.1, 1.5 and 1.13. Let  $\kappa$  be an arbitrarily small positive constant and  $\varepsilon_2$  be any sufficiently small positive constant. There is  $N_0 = N_0(n)$ , such that for all  $N \geq N_0$  and all  $z \in \mathbf{D}(N, \kappa, \varepsilon_2)$ , we have

$$\mathbb{E}|G_{ab}(z)|^{2n} \leq N^{C_0} \left( \delta_{ab} + \frac{1}{(N\eta)^n} \right), \quad \forall a, b = 1, \dots, N \quad (2.1)$$

for some positive constant  $C_0$  uniform in  $n$  and  $z$ .

By the definition of stochastic domination in Definition 1.6, we can get the following corollary from Lemma 2.1 immediately.

**Corollary 2.2.** Under the assumptions of Lemma 2.1, we have

$$|G_{ab}(z)| \prec \delta_{ab} + \frac{1}{\sqrt{N}\eta}, \quad \forall a, b = 1, \dots, N \quad (2.2)$$

uniformly on  $\mathbf{D}(N, \kappa, \varepsilon_2)$ .

In the sequel, at first, we prove Lemma 2.1 from Theorem 1.14 via the Green's function comparison strategy. Then we prove Theorem 1.9, using Lemma 2.1. Finally, we will show that Theorem 1.11 follows from Theorem 1.9 simply.

**2.1. Green's function comparison: Proof of Lemma 2.1.** To show (2.1), we use Lindeberg's replacement strategy to compare the Green's functions of the Gaussian case and the general case. That means, we will replace the entries of  $H^g$  by those of  $H$  one by one, and compare the Green's functions step by step. Choose and fix a bijective ordering map

$$\varpi : \{(i, j) : 1 \leq i \leq j \leq N\} \rightarrow \{1, \dots, \varsigma(N)\}, \quad \varsigma(N) := \frac{N(N+1)}{2}. \quad (2.3)$$

Then we use  $H_k$  to represent the  $N \times N$  random Hermitian matrix whose  $(i, j)$ -th entry is  $h_{ij}$  if  $\varpi(i, j) \leq k$ , and is  $h_{ij}^g$  otherwise. Especially, we have  $H_0 = H^g$  and  $H_{\varsigma(N)} = H$ . Correspondingly, we define the Green's functions by

$$G_k(z) := \left( H_k - z \right)^{-1}, \quad k = 1, \dots, \varsigma(N).$$

Fix  $k$  and denote

$$\varpi^{-1}(k) = (a, b). \quad (2.4)$$

Then, we write

$$\begin{aligned} H_{k-1} &= H_k^0 + \mathbf{V}_{ab}, & \mathbf{V}_{ab} &:= \left( 1 - \frac{\delta_{ab}}{2} \right) (h_{ab}^g \mathbf{e}_a \mathbf{e}_b^* + h_{ba}^g \mathbf{e}_b \mathbf{e}_a^*), \\ H_k &= H_k^0 + \mathbf{W}_{ab}, & \mathbf{W}_{ab} &:= \left( 1 - \frac{\delta_{ab}}{2} \right) (h_{ab} \mathbf{e}_a \mathbf{e}_b^* + h_{ba} \mathbf{e}_b \mathbf{e}_a^*), \end{aligned}$$

where  $H_k^0$  is obtained via replacing  $h_{ab}$  and  $h_{ba}$  by 0 in  $H_k$  (or replacing  $h_{ab}^g$  and  $h_{ba}^g$  by 0 in  $H_{k-1}$ ). In addition, we denote

$$G_k^0(z) = (H_k^0 - z)^{-1}.$$

Set  $\varepsilon_3 \equiv \varepsilon_3(\gamma, \varepsilon_1)$  to be a sufficiently small positive constant, satisfying (say)

$$\varepsilon_3 \leq \frac{1}{100} \cdot \frac{\varepsilon_1}{5 + 2\gamma + \varepsilon_1}, \quad (2.5)$$

where  $\gamma$  is from Assumption 1.1 (iii) and  $\varepsilon_1$  is from (1.7). For simplicity, we introduce the following parameters for  $\ell = 1, \dots, \varsigma(N)$  and  $i, j = 1, \dots, N$ ,

$$\widehat{\Theta}_0 := N^{C_0}, \quad \widehat{\Theta}_{\ell,ij} := \widehat{\Theta}_0 \left( 1 + C \left( \frac{N^{\varepsilon_3}}{\sqrt{M}} \right)^5 \right)^\ell \prod_{\varpi(a,b) \leq \ell} \left( 1 + C \delta_{\{i,j\}\{a,b\}} \left( \frac{N^{\varepsilon_3} \sqrt{N\eta}}{\sqrt{M}} \right)^5 \right), \quad (2.6)$$

where  $C$  is a positive constant. Here we used the notation  $\delta_{IJ} = 1$  if two index sets  $I$  and  $J$  are the same and  $\delta_{IJ} = 0$  otherwise. It is easy to see that for  $\eta \leq M^{-1}N^{\varepsilon_2}$ , we have

$$\widehat{\Theta}_{\ell,ij} \leq 2\widehat{\Theta}_0, \quad \forall \ell = 1, \dots, \varsigma(N), \quad i, j = 1, \dots, N, \quad (2.7)$$

by using (1.8).

Now, we compare  $G_{k-1}(z)$  and  $G_k(z)$ . We will prove the following lemma.

**Lemma 2.3.** *Suppose that the assumptions in Lemma 2.1 hold. Additionally, we assume that for some sufficiently small positive constant  $\varepsilon_3$  satisfying (2.5),*

$$|(G_\ell)_{ij}(z)| \prec N^{\varepsilon_3}, \quad |(G_\ell^0)_{ij}(z)| \prec N^{\varepsilon_3}, \quad \forall \ell = 1, \dots, \varsigma(N), \quad \forall i, j = 1, \dots, N \quad (2.8)$$

uniformly for  $z \in \mathbf{D}(N, \kappa, \varepsilon_2)$ . Let  $n \in \mathbb{N}$  be any given integer. Then, if

$$\mathbb{E}|(G_{k-1})_{ij}(z)|^{2n} \leq \widehat{\Theta}_{k-1,ij} \left( \delta_{ij} + \frac{1}{(N\eta)^n} \right), \quad \forall i, j = 1, \dots, N, \quad (2.9)$$

we also have

$$\mathbb{E}|(G_k)_{ij}(z)|^{2n} \leq \widehat{\Theta}_{k,ij} \left( \delta_{ij} + \frac{1}{(N\eta)^n} \right), \quad \forall i, j = 1, \dots, N \quad (2.10)$$

for any  $k = 1, \dots, \varsigma(N)$ .

*Proof of Lemma 2.3.* Fix  $k$  and omit the argument  $z$  from now on. At first, under the conditions (2.8) and (2.9), we show that

$$\mathbb{E}|(G_k^0)_{ij}|^{2n} \leq 3\widehat{\Theta}_0 \left( \delta_{ij} + \frac{1}{(N\eta)^n} \right), \quad \forall i, j = 1, \dots, N. \quad (2.11)$$

To see this, we use the expansion with (2.4)

$$(G_k^0)_{ij} = (G_{k-1})_{ij} + (G_{k-1}W_{ab}G_k^0)_{ij},$$

which implies that for a sufficiently large constant  $D > 0$

$$\mathbb{E}|(G_k^0)_{ij}|^{2n} \leq \mathbb{E} \left| (G_{k-1})_{ij} + \frac{N^{2\varepsilon_3}}{\sqrt{M}} \right|^{2n} + \eta^{-2n} N^{-D} \leq 3\widehat{\Theta}_0 \left( \delta_{ij} + \frac{1}{(N\eta)^n} \right),$$

where the first step follows from (1.12), (2.8), Definition 1.6 and the trivial bound  $\eta^{-1}$  for the Green's functions, and the second step follows from (2.9), (2.7) and the fact  $N^{2\varepsilon_3}/\sqrt{M} \ll 1/\sqrt{N\eta}$  for  $z \in \mathbf{D}(N, \kappa, \varepsilon_2)$ .

Now, recall (2.4) again and expand  $G_{k-1}(z)$  and  $G_k(z)$  around  $G_k^0(z)$ , namely

$$\begin{aligned} G_{k-1} &= G_k^0 + \sum_{\ell=1}^m (-1)^\ell (G_k^0 V_{ab})^\ell G_k^0 + (-1)^{m+1} (G_k^0 V_{ab})^{m+1} G_{k-1}, \\ G_k &= G_k^0 + \sum_{\ell=1}^m (-1)^\ell (G_k^0 W_{ab})^\ell G_k^0 + (-1)^{m+1} (G_k^0 W_{ab})^{m+1} G_k. \end{aligned} \quad (2.12)$$

We always choose  $m$  to be sufficiently large, depending on  $\varepsilon_3$  but independent of  $N$ . Then, we can write

$$(G_{k-1})_{ij} = (G_k^0)_{ij} + \sum_{\ell=1}^m R_{\ell,ij} + \tilde{R}_{m+1,ij},$$

$$(G_k)_{ij} = (G_k^0)_{ij} + \sum_{\ell=1}^m S_{\ell,ij} + \tilde{S}_{m+1,ij}, \quad (2.13)$$

where

$$\begin{aligned} R_{\ell,ij} &:= (-1)^\ell \left( (G_k^0 V_{ab})^\ell G_k^0 \right)_{ij}, \quad S_{\ell,ij} := (-1)^\ell \left( (G_k^0 W_{ab})^\ell G_k^0 \right)_{ij}, \quad \ell = 1, \dots, m, \\ \tilde{R}_{m+1,ij} &:= (-1)^{m+1} \left( (G_k^0 V_{ab})^{m+1} G_{k-1} \right)_{ij}, \quad \tilde{S}_{m+1,ij} := (-1)^{m+1} \left( (G_k^0 W_{ab})^{m+1} G_k \right)_{ij}. \end{aligned} \quad (2.14)$$

At first, by taking  $m$  sufficiently large, from (2.8) and (1.12), we have the trivial bound

$$|\tilde{R}_{m+1,ij}|, |\tilde{S}_{m+1,ij}| \prec M^{-\frac{m+1}{2}} N^{(m+2)\varepsilon_3} \ll \frac{1}{M^3 \sqrt{N\eta}}. \quad (2.15)$$

For  $R_{\ell,ij}$  and  $S_{\ell,ij}$ , we split the discussion into off-diagonal case and diagonal case. In the case of  $i \neq j$ , we keep the first and the last factors of the terms in the expansions of  $((G_k^0 V_{ab})^\ell G_k^0)_{ij}$  and  $((G_k^0 W_{ab})^\ell G_k^0)_{ij}$ , namely,  $(G_k^0)_{ij'}$  and  $(G_k^0)_{i'j}$  for some  $i', j' = a$  or  $b$ , and bound the factors in between by using (1.12) and (2.8), resulting the bound

$$|R_{\ell,ij}|, |S_{\ell,ij}| \prec M^{-\frac{\ell}{2}} N^{(\ell-1)\varepsilon_3} \sum_{i', j'=a,b} |(G_k^0)_{ij'} (G_k^0)_{i'j}|, \quad \ell = 1, \dots, m. \quad (2.16)$$

For  $i = j$ , we only keep the first factor of the terms in the expansions of  $((G_k^0 V_{ab})^\ell G_k^0)_{ii}$  and  $((G_k^0 W_{ab})^\ell G_k^0)_{ii}$ , and bound the others by using (1.12) and (2.8), resulting the bound

$$|R_{\ell,ii}|, |S_{\ell,ii}| \prec M^{-\frac{\ell}{2}} N^{\ell\varepsilon_3} \left( |(G_k^0)_{ia}| + |(G_k^0)_{ib}| \right), \quad \ell = 1, \dots, m. \quad (2.17)$$

Observe that, in case  $i \neq j$ , if  $\{i, j\} \neq \{a, b\}$ , at least one of  $(G_k^0)_{ij'}$  and  $(G_k^0)_{i'j}$  is an off-diagonal entry of  $G_k^0$  for  $i', j' = a$  or  $b$ .

Now we compare the  $2n$ -th moment of  $|(G_{k-1})_{ij}|$  and  $|(G_k)_{ij}|$ . At first, we write

$$\mathbb{E}|(G_d)_{ij}|^{2n} = \mathbb{E}((G_d)_{ij})^n \overline{((G_d)_{ij})}^n, \quad d = k-1, k \quad (2.18)$$

By substituting the expansion (2.13) into (2.18), we can write

$$\mathbb{E}|(G_d)_{ij}|^{2n} = \mathbf{A}(i, j) + \mathbf{R}_d(i, j), \quad d = k-1, k, \quad (2.19)$$

where  $\mathbf{A}(i, j)$  is the sum of the terms which depend only on  $H_k^0$  and the first four moments of  $h_{ab}$ , and  $\mathbf{R}_d(i, j)$  is the sum of all the other terms. We claim that  $\mathbf{R}_d(i, j)$  satisfies the bound

$$|\mathbf{R}_d(i, j)| \leq C \hat{\Theta}_0 \left( \frac{N^{\varepsilon_3}}{\sqrt{M}} \right)^5 \left( \delta_{ij} + \frac{\delta_{\{i,j\}\{a,b\}}}{(N\eta)^{n-\frac{5}{2}}} + \frac{1}{(N\eta)^n} \right), \quad d = k-1, k, \quad (2.20)$$

for some positive constant  $C$ . Now, we verify (2.20). According to (2.11) and the fact that the sequence  $R_{1,ij}, \dots, R_{m,ij}, \tilde{R}_{m+1,ij}$ , as well as  $S_{1,ij}, \dots, S_{m,ij}, \tilde{S}_{m+1,ij}$ , decreases by a factor  $N^{\varepsilon_3}/\sqrt{M}$  in magnitude, it is not difficult to check the leading order terms of  $\mathbf{R}_{k-1}(i, j)$  are of the form

$$\mathbb{E} \left( (G_k^0)_{ij} \right)^p \overline{\left( (G_k^0)_{ij} \right)}^{2n-p-\sum_{\ell=1}^5 (q_\ell + q'_\ell)} \prod_{\ell=1}^5 R_{\ell,ij}^{q_\ell} \bar{R}_{\ell,ij}^{q'_\ell}, \quad (2.21)$$

and those of  $\mathbf{R}_k(i, j)$  are of the form

$$\mathbb{E} \left( (G_k^0)_{ij} \right)^p \overline{\left( (G_k^0)_{ij} \right)}^{2n-p-\sum_{\ell=1}^5 (q_\ell + q'_\ell)} \prod_{\ell=1}^5 S_{\ell,ij}^{q_\ell} \bar{S}_{\ell,ij}^{q'_\ell}, \quad (2.22)$$

with some  $p, q_\ell, q'_\ell \in \mathbb{N}$  such that

$$\sum_{\ell=1}^5 \ell(q_\ell + q'_\ell) = 5, \quad 0 \leq p \leq 2n - \sum_{\ell=1}^5 (q_\ell + q'_\ell). \quad (2.23)$$

Every other term has at least 6 factors of  $h_{ab}$  or  $h_{ab}^g$  or their conjugates, thus their sizes are typically controlled by  $M^{-3}(N\eta)^{-n}$ , i.e. they are subleading. Hence, it suffices to bound (2.21) and (2.22). In the sequel, we only estimate (2.21) in details, (2.22) can be handled in the same manner.

Now, the five factors of  $h_{ab}$  or  $h_{ba}$  within the  $R_{\ell,ij}$ 's in (2.21) are independent of the rest and estimated by  $M^{-5/2}$ . For the remaining factors from  $G_k^0$ , we use (2.11) to bound  $2n$  of them and use (2.8) to bound the rest. In the case that  $i \neq j$  and  $\{i, j\} \neq \{a, b\}$ , by the discussion above, we must have an off-diagonal entry of  $G_k^0$  in the product  $(G_k^0)_{ij'}(G_k^0)_{i'j}$  for any choice of  $i', j' = a$  or  $b$ . Then, in the bound for  $R_{\ell,ij}$  in (2.16), for each  $(G_k^0)_{ij'}(G_k^0)_{i'j}$ , we keep the off-diagonal entry and bound the other by  $N^{\varepsilon_3}$  from assumption (2.8). Hence, by using (2.16) and (2.23), we see that for some  $i_r, j_r \in \{i, j, a, b\}$  with  $i_r \neq j_r$ ,  $r = 1, \dots, \sum(q_\ell + q'_\ell)$ , the following bound holds

$$(2.21) \leq \left(\frac{N^{\varepsilon_3}}{\sqrt{M}}\right)^5 \mathbb{E} \left( |(G_k^0)_{ij}|^{2n - \sum_{\ell=1}^5 (q_\ell + q'_\ell)} \prod_{r=1}^{\sum_{\ell=1}^5 (q_\ell + q'_\ell)} |(G_k^0)_{i_r j_r}| \right) \leq 3 \left(\frac{N^{\varepsilon_3}}{\sqrt{M}}\right)^5 \frac{\widehat{\Theta}_0}{(N\eta)^n}, \quad (2.24)$$

where the last step follows from (2.11) and Hölder's inequality. In case of  $i \neq j$  but  $\{i, j\} = \{a, b\}$ , we keep an entry in the product  $(G_k^0)_{ij'}(G_k^0)_{i'j}$  and bound the other by  $N^{\varepsilon_3}$ . We remark here in this case the entry being kept can be either diagonal or off-diagonal. Consequently, for some  $i_r, j_r \in \{i, j, a, b\}$ ,  $r = 1, \dots, \sum(q_\ell + q'_\ell)$ , we have the bound

$$(2.21) \leq \left(\frac{N^{\varepsilon_3}}{\sqrt{M}}\right)^5 \mathbb{E} \left( |(G_k^0)_{ij}|^{2n - \sum_{\ell=1}^5 (q_\ell + q'_\ell)} \prod_{r=1}^{\sum_{\ell=1}^5 (q_\ell + q'_\ell)} |(G_k^0)_{i_r j_r}| \right) \leq 3 \left(\frac{N^{\varepsilon_3}}{\sqrt{M}}\right)^5 \frac{\widehat{\Theta}_0}{(N\eta)^{n - \frac{5}{2}}} \quad (2.25)$$

by using (2.11) and Hölder's inequality again. Hence, we have shown (2.20) in the case of  $i \neq j$ . For  $i = j$ , it is analogous to show

$$(2.21) \leq 3 \left(\frac{N^{\varepsilon_3}}{\sqrt{M}}\right)^5 \widehat{\Theta}_0 \quad (2.26)$$

by using (2.11), (2.17) and Hölder's inequality. Hence, we verified (2.20). Consequently, by Assumption 1.5, (2.19) and (2.20) we have

$$\left| \mathbb{E} |(G_{k-1})_{ij}|^{2n} - \mathbb{E} |(G_k)_{ij}|^{2n} \right| \leq C \widehat{\Theta}_0 \left(\frac{N^{\varepsilon_3}}{\sqrt{M}}\right)^5 \left( \delta_{ij} + \frac{\delta_{\{i,j\}\{a,b\}}}{(N\eta)^{n - \frac{5}{2}}} + \frac{1}{(N\eta)^n} \right),$$

which together with the assumption (2.9) for  $\mathbb{E} |(G_{k-1})_{ij}|^{2n}$  and the definition of  $\widehat{\Theta}_{\ell,ij}$ 's in (2.6), we can get

$$\mathbb{E} |(G_k)_{ij}|^{2n} \leq \widehat{\Theta}_{k,ij} \left( \delta_{ij} + \frac{1}{(N\eta)^n} \right).$$

Hence, we completed the proof of Lemma 2.3.  $\square$

To show (2.1), we also need the following lemma.

**Lemma 2.4.** *Suppose that the assumptions in Lemma 2.1 hold. Fix the indices  $a, b \in \{1, \dots, N\}$ . Let  $H^0$  be a matrix obtained from  $H$  with its  $(a, b)$ -th entry replaced by 0. Then, if for some  $\eta_0 \geq 1/N$  there exists*

$$|G_{ii}(z)| \prec 1, \quad |(G^0)_{ii}(z)| \prec 1 \quad \text{for } \eta \geq \eta_0, \quad \forall i = 1, \dots, N, \quad (2.27)$$

then we also have

$$|G_{ij}(z)| \prec \frac{\eta_0}{\eta}, \quad |(G^0)_{ij}(z)| \prec \frac{\eta_0}{\eta}, \quad \text{for } \frac{1}{N} < \eta \leq \eta_0, \quad \forall i, j = 1, \dots, N.$$

*Proof of Lemma 2.4.* The proof is almost the same as the discussion on pages 2311–2312 in [10]. For the convenience of the reader, we sketch it below. At first, according to the discussion below (4.28) in [10], for any  $i, j = 1, \dots, N$ , we have

$$|G_{ij}(E + \mathbf{i}\eta)| \leq C \max_{\ell} \sum_{k \geq 0} \operatorname{Im} G_{\ell\ell}(E + \mathbf{i}2^k\eta).$$

Now, we set

$$k_1 := \max\{k : 2^k\eta < \eta_0\}, \quad k_2 := \max\{k : 2^k\eta < 1\}.$$

According to our assumption, both  $k_1$  and  $k_2$  are of the order  $\log N$ . Now, we have

$$\begin{aligned} \sum_{k \geq 0} \operatorname{Im} G_{\ell\ell}(E + \mathbf{i}2^k\eta) &= \sum_{k=0}^{k_1} \operatorname{Im} G_{\ell\ell}(E + \mathbf{i}2^k\eta) + \sum_{k=k_1}^{k_2} \operatorname{Im} G_{\ell\ell}(E + \mathbf{i}2^k\eta) + \sum_{k=k_2+1}^{\infty} \operatorname{Im} G_{\ell\ell}(E + \mathbf{i}2^k\eta) \\ &\prec \frac{\eta_0}{\eta} \sum_{k=0}^{k_1} \frac{1}{2^k} \operatorname{Im} G_{\ell\ell}(E + \mathbf{i}\eta_0) + (k_2 - k_1) + 1 \prec \frac{\eta_0}{\eta} \end{aligned}$$

where in the second step, we used the fact that the function  $y \mapsto y \operatorname{Im} G_{\ell\ell}(E + \mathbf{i}y)$  is monotonically increasing, the condition (2.27) and the fact  $\eta \leq \eta_0$ . Hence, we conclude the proof of Lemma 2.4.  $\square$

Now, with Theorem 1.14, Lemma 2.3 and Lemma 2.4, we can prove Lemma 2.1.

*Proof for Lemma 2.1.* The proof relies on the following bootstrap argument, namely, we show that once

$$|(G_{\ell})_{ij}| \prec 1, \quad \forall \ell = 1, \dots, \varsigma(N), \quad \forall i, j = 1, \dots, N \quad (2.28)$$

holds for  $\eta \geq \eta_0$  with  $\eta_0 \in [N^{-1+\varepsilon_2+\varepsilon_3}, M^{-1}N^{\varepsilon_2}]$ , it also holds for  $\eta \geq \eta_0 N^{-\varepsilon_3}$  for any  $\varepsilon_3$  satisfying (2.5). Assuming (2.28) holds for  $\eta \geq \eta_0$ , we see that

$$\begin{aligned} \max_{i,j} |(G_{\ell}^0)_{ij}| &= \max_{i,j} |(G_{\ell})_{ij} + ((G_{\ell})\mathbf{W}_{ab}G_{\ell}^0)_{ij}| \\ &\prec \max_{i,j} |(G_{\ell})_{ij}| + \frac{1}{\sqrt{M}} \max_{i,j} |(G_{\ell})_{ij}| \cdot \max_{i,j} |(G_{\ell}^0)_{ij}| \prec 1 + \frac{1}{\sqrt{M}} \max_{i,j} |(G_{\ell}^0)_{ij}|. \end{aligned}$$

Consequently, for  $\eta \geq \eta_0$ , we also have

$$|(G_{\ell}^0)_{ij}| \prec 1, \quad \forall \ell = 1, \dots, \varsigma(N), \quad \forall i, j = 1, \dots, N. \quad (2.29)$$

Therefore, (2.27) holds. Then, by Lemma 2.4, we see that (2.8) holds for  $\eta \geq \eta_0 N^{-\varepsilon_3}$ . Furthermore, by Lemma 2.3 and Theorem 1.14 for  $G_0$ , i.e. the Gaussian case, one can get that for any given  $n$ ,

$$\begin{aligned} \mathbb{E}|(G_{\ell})_{ij}|^{2n} &\leq \widehat{\Theta}_{\ell,ij} \left( \delta_{ij} + \frac{1}{(N\eta)^n} \right) \leq 2\widehat{\Theta}_0 \left( \delta_{ij} + \frac{1}{(N\eta)^n} \right), \quad \text{for } M^{-1}N^{\varepsilon_2} \geq \eta \geq \eta_0 N^{-\varepsilon_3}, \\ &\quad \forall \ell = 1, \dots, \varsigma(N), \quad \forall i, j = 1, \dots, N. \end{aligned} \quad (2.30)$$

Note that since (2.30) holds for any given  $n$ , we get (2.28) for  $M^{-1}N^{\varepsilon_2} \geq \eta \geq \eta_0 N^{-\varepsilon_3}$ .

Now we start from  $\eta_0 = M^{-1}N^{\varepsilon_2}$ . By Proposition 1.7 we see that (2.28) holds for all  $\eta \geq \eta_0$ . Then we can use the bootstrap argument above finitely many times to show (2.28) holds for all  $\eta \geq N^{-1+\varepsilon_2}$ . Consequently, we have (2.8) for all  $\eta \geq N^{-1+\varepsilon_2}$ . Then, Lemma 2.1 follows from Lemma 2.3 and Theorem 1.14 immediately.  $\square$



**2.2. Proof of Theorem 1.9.** Without loss of generality we can assume that  $M \leq N(\log N)^{-10}$ , otherwise, Proposition 1.7 implies (1.16) immediately. Now, recalling the notation defined in (1.3), we denote the Green's function of  $H^{(i)}$  as

$$G^{(i)}(z) := (H^{(i)} - z)^{-1},$$

with a little abuse of notation. We only need to consider the diagonal entries  $G_{ii}$  below, since the bound for the off-diagonal entries of  $G(z)$  is implied by (2.1) directly. Set

$$\Delta_i \equiv \Delta_i(z) := \frac{1}{G_{ii}} + z + \sum_a \sigma_{ai}^2 G_{aa}. \quad (2.31)$$

We introduce the notation

$$\Lambda_d \equiv \Lambda_d(z) := \max_i |G_{ii}(z) - m_{sc}(z)|.$$

We have the following lemma.

**Lemma 2.5.** *Suppose that  $H$  satisfies Assumptions 1.1, 1.5 and 1.13. We have*

$$|\Delta_i(z)| \prec \frac{1}{\sqrt{N\eta}}, \quad \forall i = 1, \dots, N, \quad (2.32)$$

uniformly for  $z \in \mathbf{D}(N, \kappa, \varepsilon_2)$ .

The proof of Lemma 2.5 will be postponed. Using Lemma 2.5, we see that, with high probability, (2.31) is a small perturbation of the self-consistent equation of  $m_{sc}$ , i.e. (1.13), considering  $\sum_a \sigma_{ai}^2 = 1$ . To control  $\Lambda_d$ , we use a *continuity argument* from [12].

We remind here that in the sequel, the parameter set of the stochastic dominance is always  $\mathbf{D}(N, \kappa, \varepsilon_2)$ , without further mention. We need to show that

$$\Lambda_d \prec \frac{1}{\sqrt{N\eta}}, \quad (2.33)$$

and first we claim that it suffices to show that

$$\mathbf{1}(\Lambda_d \leq N^{-\frac{\varepsilon_2}{4}}) \Lambda_d \prec \frac{1}{\sqrt{N\eta}}. \quad (2.34)$$

Indeed, if (2.34) were proven, we see that with high probability either  $\Lambda_d > N^{-\frac{\varepsilon_2}{4}}$  or  $\Lambda_d \prec 1/\sqrt{N\eta} \leq N^{-\frac{\varepsilon_2}{2}}$  for  $z \in \mathbf{D}(N, \kappa, \varepsilon_2)$ . That means, there is a gap in the possible range of  $\Lambda_d$ . Now, choosing  $\varepsilon$  in (1.14) to be sufficiently small, we are able to get for  $\eta = M^{-1}N^{\varepsilon_2}$ ,

$$\Lambda_d \prec N^{-\frac{\varepsilon_2}{2}}, \quad \forall E \in [-2 + \kappa, 2 - \kappa], \quad \forall i = 1, \dots, N. \quad (2.35)$$

By the fact that  $\Lambda_d$  is continuous in  $z$ , we see that with high probability,  $\Lambda_d$  can only stay in one side of the range, namely, (2.33) holds. The rigorous details of this argument involve considering a fine discrete grid of the  $z$ -parameter and using that  $G(z)$  is Lipschitz continuous (albeit with a large Lipschitz constant  $1/\eta$ ). The details are found in Section 5.3 of [12].

Hence, what remains is to verify (2.34). The proof of (2.34) is almost the same as that for Lemma 3.5 in [14]. For the convenience of the reader, we sketch it below without reproducing the details. We set

$$\bar{m} \equiv \bar{m}(z) := \frac{1}{N} \sum_{i=1}^N G_{ii}(z), \quad \mathbf{u}_i \equiv \mathbf{u}_i(z) := G_{ii} - \bar{m}, \quad i = 1, \dots, N.$$

We also denote  $\vec{\mathbf{u}} := (\mathbf{u}_1, \dots, \mathbf{u}_N)$ . By the assumption  $\Lambda_d \leq N^{-\frac{\varepsilon_2}{4}}$ , we have

$$\mathbf{u}_i = O(N^{-\frac{\varepsilon_2}{4}}). \quad (2.36)$$

Now we rewrite (2.31) as

$$0 = G_{ii} + \frac{1}{z + \sum_a \sigma_{ai}^2 G_{aa} - \Delta_i} =: G_{ii} + \frac{1}{z + \bar{m}(z)} + \Omega_i. \quad (2.37)$$

By using (2.32), Lemma 5.1 in [14], and the assumption  $\Lambda_d \leq N^{-\frac{\varepsilon_2}{4}}$ , we can show that

$$\Omega_i = -\frac{\sum_a \sigma_{ai}^2 u_a}{(z + \bar{m}(z))^2} + O(\|\bar{\mathbf{u}}\|_\infty^2) + O(\max_i |\Delta_i|). \quad (2.38)$$

One can refer to the derivation of (5.14) in [14] for more details. Averaging over  $i$  for (2.37) and (2.38) leads to

$$\bar{m}(z) + \frac{1}{z + \bar{m}(z)} = -\Omega. \quad (2.39)$$

and

$$\Omega := \frac{1}{N} \sum_{i=1}^N \Omega_i = O(\|\bar{\mathbf{u}}\|_\infty^2) + O(\max_i |\Delta_i|) \quad (2.40)$$

Plugging (2.36) and (2.32) into (2.40) yields

$$|\Omega| \prec N^{-\frac{\varepsilon_2}{2}}. \quad (2.41)$$

Using (2.41), the fact  $|\bar{m}(z) - m_{sc}(z)| \leq \Lambda_d \leq N^{-\frac{\varepsilon_2}{4}}$ , and Lemma 5.2 in [14], to (2.39), we have

$$|\bar{m}(z) - m_{sc}(z)| \leq |\Omega| = O(\|\bar{\mathbf{u}}\|_\infty^2) + O(\max_i |\Delta_i|), \quad (2.42)$$

where in the first step we have used the fact that  $z \in \mathbf{D}(N, \kappa, \varepsilon_2)$  thus away from the edges of the semicircle law. Now, we combine (2.37), (2.38) and (2.39), resulting

$$u_i = \frac{\sum_a \sigma_{ai}^2 u_a}{(z + \bar{m}(z))^2} + \Omega + O(\|\bar{\mathbf{u}}\|_\infty^2) + O(\max_i |\Delta_i|) = w_i + \frac{\sum_a \sigma_{ai}^2 u_a}{(z + m_{sc}(z))^2}, \quad i = 1, \dots, N. \quad (2.43)$$

We just take the above identity as the definition of  $w_i$ . Analogously, we set  $\vec{w} := (w_1, \dots, w_N)'$ . Then (2.40) and (2.43) imply

$$\begin{aligned} \|\vec{w}\|_\infty &= O(\|\bar{\mathbf{u}}\|_\infty^2) + O(\max_i |\Delta_i|) + O(\|\bar{\mathbf{u}}\|_\infty \cdot |(z + \bar{m}(z))^{-2} - (z + m_{sc}(z))^{-2}|) \\ &\leq O(\|\bar{\mathbf{u}}\|_\infty^2) + O(\max_i |\Delta_i|) + O(\|\bar{\mathbf{u}}\|_\infty \cdot |\bar{m}(z) - m_{sc}(z)|) \\ &\leq O(\|\bar{\mathbf{u}}\|_\infty^2) + O(\max_i |\Delta_i|) \end{aligned} \quad (2.44)$$

where the second step follows from the fact  $|z + m_{sc}(z)| \geq 1$  in  $\mathbf{D}(N, \kappa, \varepsilon_2)$  (see (5.1) in [14] for instance), (2.41) and (2.42), and in the last step we used (2.42) again.

Now, using the fact  $m_{sc}^2(z) = (m_{sc}(z) + z)^{-2}$  (see (1.13)), we rewrite (2.43) in terms of the matrix  $\mathcal{T}$  introduced in (1.5) as

$$\vec{u} = (1 - m_{sc}^2(z)\mathcal{T})^{-1} \vec{w}.$$

Consequently, we have

$$\|\vec{u}\|_\infty \leq \left\| (1 - m_{sc}^2(z)\mathcal{T})^{-1} \right\|_{\ell^\infty \rightarrow \ell^\infty} \|\vec{w}\|_\infty := \Gamma(z) \|\vec{w}\|_\infty. \quad (2.45)$$

Then for  $z \in \mathbf{D}(N, \kappa, \varepsilon_2)$ , using (1.6) and Proposition A.2 (ii) in [12] (with  $\delta_- = 1$  and  $\theta > c$ ), we can get

$$\Gamma(z) = O(\log N). \quad (2.46)$$

Plugging (2.46) and (2.44) into (2.45) yields

$$\|\vec{u}\|_\infty \prec O(\|\bar{\mathbf{u}}\|_\infty^2 + O(\max_i |\Delta_i|)) \prec \|\bar{\mathbf{u}}\|_\infty^2 + \frac{1}{\sqrt{N\eta}},$$

where the second step follows from (2.32). Then (2.36) further implies that

$$\|\vec{u}\|_\infty \prec \frac{1}{\sqrt{N\eta}},$$

which together with (2.42) and (2.32) also implies

$$|\bar{m}(z) - m_{sc}(z)| \prec \frac{1}{\sqrt{N\eta}}.$$

Hence

$$\Lambda_d \leq \|\vec{u}\|_\infty + |\bar{m}(z) - m_{sc}(z)| \prec \frac{1}{\sqrt{N\eta}}.$$

Therefore, we completed the proof of Theorem 1.9.

*Proof of Lemma 2.5.* For simplicity, we omit the variable  $z$  from the notation below. At first, we recall the elementary identity by Schur's complement, namely,

$$G_{ii} = \frac{1}{h_{ii} - z - (\mathbf{h}_i^{(i)})^* G^{(i)} \mathbf{h}_i^{(i)}}. \quad (2.47)$$

where we used the notation  $\mathbf{h}_i^{(i)}$  to denote the  $i$ -th column of  $H$ , with the  $i$ -th component deleted. Now, we use the identity for  $a, b \neq i$  (see Lemma 4.5 in [12] for instance),

$$G_{ab}^{(i)} = G_{ab} - G_{ai} G_{ib} (G_{ii})^{-1} = G_{ab} - G_{ai} G_{ib} \left( h_{ii} - z - (\mathbf{h}_i^{(i)})^* G^{(i)} \mathbf{h}_i^{(i)} \right). \quad (2.48)$$

By using (1.10) and the large deviation estimate for the quadratic form (see Theorem C.1 of [12] for instance), we have

$$\left| (\mathbf{h}_i^{(i)})^* G^{(i)} \mathbf{h}_i^{(i)} - \sum_{a \neq i} \sigma_{ai}^2 \cdot G_{aa}^{(i)} \right| \prec \sqrt{\frac{1}{M} \max_a |G_{aa}^{(i)}|^2 + \max_{a \neq b} |G_{ab}^{(i)}|^2}, \quad (2.49)$$

which implies that

$$\left| (\mathbf{h}_i^{(i)})^* G^{(i)} \mathbf{h}_i^{(i)} \right| \prec \max_{a \neq i} |G_{aa}^{(i)}| + \sqrt{\frac{1}{M} \max_{a \neq i} |G_{aa}^{(i)}|^2 + \max_{a \neq b} |G_{ab}^{(i)}|^2} \leq 3 \max_{a, b \neq i} |G_{ab}^{(i)}|, \quad (2.50)$$

where we have used the fact that  $\sum_a \sigma_{ai}^2 = 1$  in the first inequality above. Plugging (1.20) and (2.50) into (2.48) and using Corollary 2.2 we obtain

$$\max_{a, b \neq i} |G_{ab}^{(i)}| \prec 1 + \frac{1}{N\eta} \left( 1 + 3 \max_{a, b \neq i} |G_{ab}^{(i)}| \right),$$

which implies

$$\max_{a, b \neq i} |G_{ab}^{(i)}| \prec 1, \quad \left| (\mathbf{h}_i^{(i)})^* G^{(i)} \mathbf{h}_i^{(i)} \right| \prec 1. \quad (2.51)$$

In addition, (1.20), (2.48) and (2.51) lead to the fact that

$$|G_{ab}(z) - G_{ab}^{(i)}(z)| \prec \frac{1}{N\eta}, \quad |G_{ab}^{(i)}(z)| \prec \delta_{ab} + \frac{1}{\sqrt{N\eta}}, \quad \forall a, b \neq i. \quad (2.52)$$

Now, using (2.31), (2.47), (2.49) and (2.52), we can see that

$$|\Delta_i| = \left| -h_{ii} + (\mathbf{h}_i^{(i)})^* G^{(i)} \mathbf{h}_i^{(i)} - \sum_a \sigma_{ai}^2 G_{aa} \right| \prec \frac{1}{\sqrt{N\eta}}. \quad (2.53)$$

Therefore, we completed the proof of Lemma 2.5.  $\square$

**2.3. Proof of Theorem 1.11.** With Theorem 1.9, we can prove Theorem 1.11 routinely. At first, due to Definition 1.6 and the fact that  $G_{ab}(z)$  and  $m_{sc}(z)$  are Lipschitz functions of  $z$  with Lipschitz constant  $\eta^{-1}$ , it is easy to strengthen (1.16) to

$$\max_{a,b} \sup_{z \in \mathbf{D}(N, \kappa, \varepsilon_2)} |G_{ab}(z) - \delta_{ab} m_{sc}(z)| \prec \frac{1}{\sqrt{N}\eta},$$

which implies that

$$\max_a \sup_{z \in \mathbf{D}(N, \kappa, \varepsilon_2)} |G_{aa}(z)| \prec C \quad (2.54)$$

for some positive constant  $C$  due to the fact that  $m_{sc}(z) \sim 1$ . Recalling the normalized eigenvector  $\mathbf{u}_i = (u_{i1}, \dots, u_{iN})$  corresponding to  $\lambda_i$ , and using the spectral decomposition, we have

$$\max_a \operatorname{Im} G_{aa}(z) = \max_a \sum_{i=1}^N \frac{|u_{ia}|^2 \eta}{|\lambda_i - E|^2 + \eta^2} = \sum_{i=1}^N \frac{\|\mathbf{u}_i\|_\infty^2 \eta}{|\lambda_i - E|^2 + \eta^2}. \quad (2.55)$$

For any  $|\lambda_i| \leq \sqrt{2} - \kappa$ , we set  $E = \lambda_i$  on the r.h.s. of (2.55) and use (2.54) to bound the l.h.s. of it. Then we obtain

$$\frac{\|\mathbf{u}_i\|_\infty^2}{\eta} \prec 1.$$

Choosing  $\eta = N^{-1+\varepsilon_2}$  above and using the fact that  $\varepsilon_2$  can be arbitrarily small, we can get (1.17). Hence, we completed the proof of Theorem 1.11.

### 3. SUPERSYMMETRIC FORMALISM AND INTEGRAL REPRESENTATION FOR THE GREEN'S FUNCTION

In this section, we will represent  $\mathbb{E}|G_{ij}(z)|^{2n}$  for the Gaussian case by a superintegral. The final representation is stated in (3.30). We make the convention here, for any real argument in an integral below, its region of the integral is always  $\mathbb{R}$ , unless specified otherwise.

**3.1. Gaussian integrals and superbosonization formulas.** Let  $\phi = (\phi_1, \dots, \phi_k)'$  be a vector of complex components,  $\psi = (\psi_1, \dots, \psi_k)'$  be a vector of Grassmann components. In addition, let  $\phi^*$  and  $\psi^*$  be the conjugate transposes of  $\phi$  and  $\psi$ , respectively. We recall the following well-known formulas for Gaussian integrals.

**Proposition 3.1** (Gaussian integrals or Wick's formulas).

- (i) Let  $A$  be a  $k \times k$  complex matrix with positive-definite Hermitian part, i.e.  $\operatorname{Re} A > 0$ . Then for any  $\ell \in \mathbb{N}$ , and  $i_1, \dots, i_\ell, j_1, \dots, j_\ell \in \{1, \dots, k\}$ , we have

$$\int \prod_{a=1}^k \frac{d\operatorname{Re} \phi_a d\operatorname{Im} \phi_a}{\pi} \exp\{-\phi^* A \phi\} \prod_{b=1}^\ell \bar{\phi}_{i_b} \phi_{j_b} = \frac{1}{\det A} \sum_{\sigma \in \mathbb{P}(\ell)} \prod_{b=1}^\ell (A^{-1})_{j_b, i_{\sigma(b)}}, \quad (3.1)$$

where  $\mathbb{P}(\ell)$  is the permutation group of degree  $\ell$ .

- (ii) Let  $B$  be any  $k \times k$  matrix. Then for any  $\ell \in \{0, \dots, k\}$ , any  $\ell$  distinct integers  $i_1, \dots, i_\ell$  and another  $\ell$  distinct integers  $j_1, \dots, j_\ell \in \{1, \dots, k\}$ , we have

$$\int \prod_{a=1}^k d\bar{\psi}_a d\psi_a \exp\{-\psi^* B \psi\} \prod_{b=1}^\ell \bar{\psi}_{i_b} \psi_{j_b} = (-1)^{\ell + \sum_{\alpha=1}^\ell (i_\alpha + j_\alpha)} \det B^{(\mathbf{l}|\mathbf{J})}, \quad (3.2)$$

where  $\mathbf{l} = \{i_1, \dots, i_\ell\}$ , and  $\mathbf{J} = \{j_1, \dots, j_\ell\}$ .

Now, we introduce the superbosonization formula for superintegrals. Let  $\chi = (\chi_{ij})$  be an  $\ell \times r$  matrix with Grassmann entries,  $\mathbf{f} = (f_{ij})$  be an  $\ell \times r$  matrix with complex entries. In addition, we denote their conjugate transposes by  $\chi^*$  and  $\mathbf{f}^*$  respectively. Let  $F$  be a function of the entries of the matrix

$$\mathcal{S}(\mathbf{f}, \mathbf{f}^*; \chi, \chi^*) := \begin{pmatrix} \chi^* \chi & \chi^* \mathbf{f} \\ \mathbf{f}^* \chi & \mathbf{f}^* \mathbf{f} \end{pmatrix}.$$

Let  $\mathcal{A}(\chi, \chi^*)$  be the Grassmann algebra generated by  $\chi_{ij}$ 's and  $\bar{\chi}_{ij}$ 's. Then we can regard  $F$  as a function defined on a complex vector space, taking values in  $\mathcal{A}(\chi, \chi^*)$ . Hence, we can and do view  $F(\mathcal{S}(\mathbf{f}, \mathbf{f}^*; \chi, \chi^*))$  as a polynomial in  $\chi_{ij}$ 's and  $\bar{\chi}_{ij}$ 's, in which the coefficients are functions of  $f_{ij}$ 's and  $\bar{f}_{ij}$ 's. Under this viewpoint, we state the assumption on  $F$  as follows.

**Assumption 3.2.** *Suppose that  $F(\mathcal{S}(\mathbf{f}, \mathbf{f}^*; \chi, \chi^*))$  is a holomorphic function of  $f_{ij}$ 's and  $\bar{f}_{ij}$ 's if they are regarded as independent variables, and  $F$  is a Schwarz function of  $\operatorname{Re} f_{ij}$ 's and  $\operatorname{Im} f_{ij}$ 's, by those we mean that all of the coefficients of  $F(\mathcal{S}(\mathbf{f}, \mathbf{f}^*; \chi, \chi^*))$ , as functions of  $f_{ij}$ 's and  $\bar{f}_{ij}$ 's, possess the above properties.*

**Proposition 3.3** (Superbosonization formula for the nonsingular case, [17]). *Suppose that  $F$  satisfies Assumption 3.2. For  $\ell \geq r$ , we have*

$$\int F \begin{pmatrix} \chi^* \chi & \chi^* \mathbf{f} \\ \mathbf{f}^* \chi & \mathbf{f}^* \mathbf{f} \end{pmatrix} d\mathbf{f} d\chi = (i\pi)^{-r(r-1)} \int d\hat{\mu}(\mathbf{x}) d\hat{\nu}(\mathbf{y}) d\omega d\xi F \begin{pmatrix} \mathbf{x} & \omega \\ \xi & \mathbf{y} \end{pmatrix} \frac{\det^\ell \mathbf{y}}{\det^\ell (\mathbf{x} - \omega \mathbf{y}^{-1} \xi)}, \quad (3.3)$$

where  $\mathbf{x} = (x_{ij})$  is a unitary matrix;  $\mathbf{y} = (y_{ij})$  is a positive-definite Hermitian matrix;  $\omega$  and  $\xi$  are two Grassmann matrices, and all of them are  $r \times r$ . Here

$$\begin{aligned} d\mathbf{f} &= \prod_{i,j} \frac{d\operatorname{Re} f_{ij} d\operatorname{Im} f_{ij}}{\pi}, & d\chi &= \prod_{i,j} d\bar{\chi}_{ij} d\chi_{ij}, \\ d\hat{\nu}(\mathbf{y}) &= \mathbf{1}(\mathbf{y} > 0) \prod_{i=1}^r dy_{ii} \prod_{j>k} d\operatorname{Re} y_{jk} d\operatorname{Im} y_{jk}, & d\omega d\xi &= \prod_{i,j=1}^r d\omega_{ij} d\xi_{ij}, \end{aligned}$$

and  $d\hat{\mu}(\cdot)$  is defined by

$$d\hat{\mu}(\mathbf{x}) = \frac{\pi^{r(r-1)/2}}{\prod_{i=1}^r i!} \cdot \prod_{i=1}^r \frac{dx_i}{2\pi i} \cdot (\Delta(x_1, \dots, x_r))^2 \cdot d\mu(V),$$

under the parametrization induced by the eigendecomposition, namely,

$$\mathbf{x} = V^* \hat{\mathbf{x}} V, \quad \hat{\mathbf{x}} = \operatorname{diag}(x_1, \dots, x_r), \quad V \in \mathring{U}(r).$$

Here  $d\mu(V)$  is the Haar measure on  $\mathring{U}(r)$ , and  $\Delta(\cdot)$  is the Vandermonde determinant. In addition, the integral w.r.p.t.  $\mathbf{x}$  ranges over  $U(2)$ , that w.r.p.t.  $\mathbf{y}$  ranges over all positive-definite matrices.

For the singular case, i.e.  $r > \ell$ , we only state the formula for the case of  $r = 2$  and  $\ell = 1$ , which is enough for our purpose. We can refer to formula (11) in [3] for the result under more general setting on  $r$  and  $\ell$ .

**Proposition 3.4** (Superbosonization formula for the singular case, [3]). *Suppose that  $F$  satisfies Assumption 3.2. If  $r = 2$  and  $\ell = 1$ , we have*

$$\int F \begin{pmatrix} \chi^* \chi & \chi^* \mathbf{f} \\ \mathbf{f}^* \chi & \mathbf{f}^* \mathbf{f} \end{pmatrix} d\mathbf{f} d\chi = \frac{-1}{\pi^2} \int d\omega d\hat{\mu}(\mathbf{x}) \cdot \mathbf{1}(y \geq 0) dy \cdot d\omega d\xi F \begin{pmatrix} \mathbf{x} & \omega \mathbf{w}^* \\ \mathbf{w} \xi & y \mathbf{w} \mathbf{w}^* \end{pmatrix} \frac{y(y - \xi \mathbf{x}^{-1} \omega)^2}{\det^2 \mathbf{x}}, \quad (3.4)$$

where  $y$  is a positive variable;  $\mathbf{x}$  is a 2-dimensional unitary matrix;  $\boldsymbol{\omega} = (\omega_1, \omega_2)'$  and  $\boldsymbol{\xi} = (\xi_1, \xi_2)$  are two vectors with Grassmann components. In addition,  $\mathbf{w}$  is a unit vector, which can be parameterized by

$$\mathbf{w} = \left( 1/\sqrt{1+|w|^2}, w/\sqrt{1+|w|^2} \right)', \quad w \in \mathbb{C}.$$

Moreover, the differentials are defined as

$$d\mathbf{w} = \frac{1}{(1+|w|^2)^2} d\text{Re}w d\text{Im}w, \quad d\boldsymbol{\omega} d\boldsymbol{\xi} = \prod_{i=1,2} d\omega_i d\xi_i.$$

In addition, the integral w.r.p.t.  $\mathbf{x}$  ranges over  $U(2)$ .

In our discussion, for  $\mathbf{w}$ , we will adopt the parametrization

$$v = 1/\sqrt{1+|w|^2}, \quad e^{i\theta} = w/|w|, \quad v \in \mathbb{I}, \quad \theta \in \mathbb{I}$$

for convenience. Accordingly, we can get

$$d\mathbf{w} = v dv d\theta.$$

**3.2. Initial representation.** For  $a = 1, 2$  and  $j = 1, \dots, W$ , we set

$$\begin{aligned} \Phi_{a,j} &= (\phi_{a,j,1}, \dots, \phi_{a,j,M})', \quad \Psi_{a,j} = (\psi_{a,j,1}, \dots, \psi_{a,j,M})' \\ \Phi_a &= (\Phi'_{a,1}, \dots, \Phi'_{a,W})', \quad \Psi_a = (\Psi'_{a,1}, \dots, \Psi'_{a,W})'. \end{aligned}$$

For each  $j$  and each  $a$ ,  $\Phi_{a,j}$  is a vector with complex components, and  $\Psi_{a,j}$  is a vector with Grassmann components. In addition, we use  $\Phi_{a,j}^*$  and  $\Psi_{a,j}^*$  to represent the conjugate transposes of  $\Phi_{a,j}$  and  $\Psi_{a,j}$  respectively. Analogously, we adopt the notation  $\Phi_a^*$  and  $\Psi_a^*$  to represent the conjugate transposes of  $\Phi_a$  and  $\Psi_a$ , respectively. We have the following integral representation for the moments of the Green's function.

**Lemma 3.5.** For any  $p, q = 1, \dots, W$  and  $\alpha, \beta = 1, \dots, M$ , we have

$$\begin{aligned} |G_{pq, \alpha\beta}(z)|^{2n} &= \frac{1}{(n!)^2} \int d\Phi d\Psi (\bar{\phi}_{1,q,\beta} \phi_{1,p,\alpha} \bar{\phi}_{2,p,\alpha} \phi_{2,q,\beta})^n \\ &\quad \times \exp \left\{ \mathbf{i} \Psi_1^* (z - H) \Psi_1 + \mathbf{i} \Phi_1^* (z - H) \Phi_1 - \mathbf{i} \Psi_2^* (\bar{z} - H) \Psi_2 - \mathbf{i} \Phi_2^* (\bar{z} - H) \Phi_2 \right\}, \end{aligned} \quad (3.5)$$

where

$$d\Phi = \prod_{a=1,2} \prod_{j=1}^W \prod_{\alpha'=1}^M \frac{d\text{Re}\phi_{a,j,\alpha'} d\text{Im}\phi_{a,j,\alpha'}}{\pi}, \quad d\Psi = \prod_{a=1,2} \prod_{j=1}^W \prod_{\alpha'=1}^M d\bar{\psi}_{a,j,\alpha'} d\psi_{a,j,\alpha'}.$$

*Proof.* By using Proposition 3.1 (i) with  $\ell = n$  and Proposition 3.1 (ii) with  $\ell = 0$ , we can get (3.5) immediately.  $\square$

**3.3. Averaging over the Gaussian random matrix.** Recall the variance profile  $\tilde{S}$  in (1.2). Now, we take expectation of the Green's function, i.e average over the random matrix. By elementary Gaussian integral, we get

$$\begin{aligned} \mathbb{E}|G_{pq, \alpha\beta}(z)|^{2n} &= \frac{1}{(n!)^2} \int d\Phi d\Psi (\bar{\phi}_{1,q,\beta} \phi_{1,p,\alpha} \bar{\phi}_{2,p,\alpha} \phi_{2,q,\beta})^n \exp \left\{ \mathbf{i} \sum_{j=1}^W (\text{Tr} \check{X}_j J Z + \text{Tr} \check{Y}_j J Z) \right\} \\ &\quad \times \exp \left\{ \frac{1}{2M} \sum_{j,k} \tilde{\mathbf{s}}_{jk} \text{Tr} \check{X}_j J \check{X}_k J - \frac{1}{2M} \sum_{j,k} \tilde{\mathbf{s}}_{jk} \text{Tr} \check{Y}_j J \check{Y}_k J \right\} \\ &\quad \times \exp \left\{ -\frac{1}{M} \sum_{j,k} \tilde{\mathbf{s}}_{jk} \text{Tr} \check{\Omega}_j J \check{\Xi}_k J \right\}, \end{aligned} \quad (3.6)$$

where

$$J = \text{diag}(1, -1), \quad Z = \text{diag}(z, \bar{z}),$$

and for each  $j = 1, \dots, W$ , the matrices  $\check{X}_j, \check{Y}_j, \check{\Omega}_j$  and  $\check{\Xi}_j$  are  $2 \times 2$  blocks of a supermatrix, namely,

$$\check{\mathcal{S}}_j = \left( \begin{array}{c|c} \check{X}_j & \check{\Omega}_j \\ \hline \check{\Xi}_j & \check{Y}_j \end{array} \right) := \left( \begin{array}{cc|cc} \Psi_{1,j}^* \Psi_{1,j} & \Psi_{1,j}^* \Psi_{2,j} & \Psi_{1,j}^* \Phi_{1,j} & \Psi_{1,j}^* \Phi_{2,j} \\ \Psi_{2,j}^* \Psi_{1,j} & \Psi_{2,j}^* \Psi_{2,j} & \Psi_{2,j}^* \Phi_{1,j} & \Psi_{2,j}^* \Phi_{2,j} \\ \hline \Phi_{1,j}^* \Psi_{1,j} & \Phi_{1,j}^* \Psi_{2,j} & \Phi_{1,j}^* \Phi_{1,j} & \Phi_{1,j}^* \Phi_{2,j} \\ \Phi_{2,j}^* \Psi_{1,j} & \Phi_{2,j}^* \Psi_{2,j} & \Phi_{2,j}^* \Phi_{1,j} & \Phi_{2,j}^* \Phi_{2,j} \end{array} \right).$$

**Remark 3.6.** The derivation of (3.6) from (3.5) is quite standard. We refer to the proof of (2.14) in [19] for more details and will not reproduce it here.

**3.4. Decomposition of the supermatrices.** From now on, we split the discussion into the following three cases

- (Case 1): Entries in the off-diagonal blocks, i.e.  $p \neq q$ ,
- (Case 2): Off-diagonal entries in the diagonal blocks, i.e.  $p = q, \quad \alpha \neq \beta$ ,
- (Case 3): Diagonal entries, i.e.  $p = q, \quad \alpha = \beta$ .

For each case, we will perform a decomposition for the supermatrix  $\check{\mathcal{S}}_j$  ( $j = p$  or  $q$ ). For a vector  $\mathbf{v}$  and some index set  $\mathbf{l}$ , we use  $\mathbf{v}^{(\mathbf{l})}$  to denote the subvector obtained by deleting the  $i$ -th component of  $\mathbf{v}$  for all  $i \in \mathbf{l}$ . Then, we adopt the notation

$$\check{\mathcal{S}}_j^{(\mathbf{l})} = \left( \begin{array}{c|c} \check{X}_j^{(\mathbf{l})} & \check{\Omega}_j^{(\mathbf{l})} \\ \hline \check{\Xi}_j^{(\mathbf{l})} & \check{Y}_j^{(\mathbf{l})} \end{array} \right), \quad \check{\mathcal{S}}_j^{[i]} = \left( \begin{array}{c|c} \check{X}_j^{[i]} & \check{\Omega}_j^{[i]} \\ \hline \check{\Xi}_j^{[i]} & \check{Y}_j^{[i]} \end{array} \right).$$

Here, for  $\mathbf{A} = \check{X}_j, \check{Y}_j, \check{\Omega}_j$  or  $\check{\Xi}_j$ , the notation  $\mathbf{A}^{(\mathbf{l})}$  is defined via replacing  $\Phi_{a,j}, \Psi_{a,j}, \Phi_{a,j}^*$  and  $\Psi_{a,j}^*$  by  $\Phi_{a,j}^{(\mathbf{l})}, \Psi_{a,j}^{(\mathbf{l})}, (\Phi_{a,j}^*)^{(\mathbf{l})}$  and  $(\Psi_{a,j}^*)^{(\mathbf{l})}$ , respectively, for  $a = 1, 2$ , in the definition of  $\mathbf{A}$ . In addition, the notation  $\mathbf{A}^{[i]}$  is defined via replacing  $\Phi_{a,j}, \Psi_{a,j}, \Phi_{a,j}^*$  and  $\Psi_{a,j}^*$  by  $\phi_{a,j,i}, \psi_{a,j,i}, \bar{\phi}_{a,j,i}$  and  $\bar{\psi}_{a,j,i}$  respectively, for  $a = 1, 2$ , in the definition of  $\mathbf{A}$ . Moreover, for  $\mathbf{A} = \check{\mathcal{S}}_j, \check{X}_j, \check{Y}_j, \check{\Omega}_j$  or  $\check{\Xi}_j$ , we will simply abbreviate  $\mathbf{A}^{\langle\{a,b\}\rangle}$  and  $\mathbf{A}^{\langle\{a\}\rangle}$  by  $\mathbf{A}^{\langle a,b \rangle}$  and  $\mathbf{A}^{\langle a \rangle}$ , respectively. Note that  $\check{\mathcal{S}}_j^{[i]}$  is of rank-one.

For Case 1, due to symmetry, we can assume  $\alpha = \beta = 1$ . Then we extract two rank-one supermatrices from  $\check{\mathcal{S}}_p$  and  $\check{\mathcal{S}}_q$  such that the quantities  $\bar{\phi}_{2,p,1}\phi_{1,p,1}$  and  $\bar{\phi}_{1,q,1}\phi_{2,q,1}$  can be expressed in terms of the entries of these supermatrices. More specifically, we decompose the supermatrices

$$\check{\mathcal{S}}_p = \check{\mathcal{S}}_p^{(1)} + \check{\mathcal{S}}_p^{[1]}, \quad \check{\mathcal{S}}_q = \check{\mathcal{S}}_q^{(1)} + \check{\mathcal{S}}_q^{[1]}. \quad (3.7)$$

Consequently, we can write

$$\bar{\phi}_{1,q,1}\phi_{1,p,1}\bar{\phi}_{2,p,1}\phi_{2,q,1} = (\check{Y}_q^{[1]})_{12}(\check{Y}_p^{[1]})_{21}. \quad (3.8)$$

For Case 2, due to symmetry, we can assume that  $\alpha = 1, \beta = 2$ . Then we extract two rank-one supermatrices from  $\check{\mathcal{S}}_p$ , namely,

$$\check{\mathcal{S}}_p = \check{\mathcal{S}}_p^{(1,2)} + \check{\mathcal{S}}_p^{[1]} + \check{\mathcal{S}}_p^{[2]}. \quad (3.9)$$

Consequently, we can write

$$\bar{\phi}_{1,p,2}\phi_{1,p,1}\bar{\phi}_{2,p,1}\phi_{2,p,2} = (\check{Y}_p^{[2]})_{12}(\check{Y}_p^{[1]})_{21}. \quad (3.10)$$

Finally, for Case 3, due to symmetry, we can assume that  $\alpha = 1$ . Then we extract only one rank-one supermatrix from  $\check{\mathcal{S}}_p$ , namely,

$$\check{\mathcal{S}}_p = \check{\mathcal{S}}_p^{(1)} + \check{\mathcal{S}}_p^{[1]}. \quad (3.11)$$



Consequently, we can write

$$\bar{\phi}_{1,p,1}\phi_{1,p,1}\bar{\phi}_{2,p,1}\phi_{2,p,1} = (\check{Y}_p^{[1]})_{12}(\check{Y}_p^{[1]})_{21} = (\check{Y}_p^{[1]})_{11}(\check{Y}_p^{[1]})_{22}.$$

Since the discussion for all three cases are similar, we will only present the details for Case 1. More specifically, in the remaining part of this section and Section 4 to Section 10, we will only treat Case 1. In Section 11, we will sum up the discussions in the previous sections and explain how to adapt them to Case 2 and Case 3, resulting a final proof of Theorem 1.14.

**3.5. Variable reduction by superbosonization formulae.** We will work with Case 1. Recall the decomposition (3.7). We use the superbosonization formulae to reduce the number of variables. We shall treat  $\check{S}_k$  ( $k \neq p, q$ ) and  $\check{S}_j^{(1)}$  ( $j = p, q$ ) on an equal footing and use the formula (3.3) with  $r = 2, \ell = M$  for the former and  $r = 2, \ell = M - 1$  for the latter, while we separate the terms  $\check{S}_j^{[1]}$  ( $j = p, q$ ) and use the formula (3.4). For simplicity, we introduce the notation

$$\tilde{S}_j = \begin{cases} \check{S}_j, & \text{if } j \neq p, q, \\ \check{S}_j^{(1)}, & \text{if } j = p, q. \end{cases} \quad (3.12)$$

Accordingly, we will use  $\tilde{X}_j$ ,  $\tilde{\Omega}_j$ ,  $\tilde{\Xi}_j$  and  $\tilde{Y}_j$  to denote four blocks of  $\tilde{S}_j$ . With this notation, we can rewrite (3.6) with  $\alpha = \beta = 1$  as

$$\begin{aligned} \mathbb{E}|G_{pq,11}(z)|^{2n} &= \frac{1}{(n!)^2} \int d\Phi d\Psi \left( \bar{\phi}_{1,q,1}\phi_{1,p,1}\bar{\phi}_{2,p,1}\phi_{2,q,1} \right)^n \exp \left\{ \mathbf{i} \sum_{j=1}^W \left( Tr \tilde{X}_j J Z + Tr \tilde{Y}_j J Z \right) \right\} \\ &\times \exp \left\{ \frac{1}{2M} \sum_{j,k} \tilde{s}_{jk} \left( Tr \tilde{X}_j J \tilde{X}_k J - Tr \tilde{Y}_j J \tilde{Y}_k J \right) \right\} \exp \left\{ -\frac{1}{M} \sum_{j,k} \tilde{s}_{jk} Tr \tilde{\Omega}_j J \tilde{\Xi}_k J \right\} \\ &\times \prod_{k=p,q} \exp \left\{ \mathbf{i} Tr \check{X}_k^{[1]} J Z + \mathbf{i} Tr \check{Y}_k^{[1]} J Z \right\} \prod_{k=p,q} \exp \left\{ \frac{1}{M} \sum_{j=1}^W \tilde{s}_{jk} \left( Tr \tilde{X}_j J \check{X}_k^{[1]} J - Tr \tilde{Y}_j J \check{Y}_k^{[1]} J \right) \right\} \\ &\times \prod_{k,\ell=p,q} \exp \left\{ \frac{\tilde{s}_{k\ell}}{2M} \left( Tr \check{X}_k^{[1]} J \check{X}_\ell^{[1]} J - Tr \check{Y}_k^{[1]} J \check{Y}_\ell^{[1]} J \right) \right\} \prod_{k,\ell=p,q} \exp \left\{ -\frac{\tilde{s}_{k\ell}}{M} Tr \check{\Omega}_k^{[1]} J \check{\Xi}_\ell^{[1]} J \right\} \\ &\times \prod_{k=p,q} \exp \left\{ -\frac{1}{M} \sum_j \tilde{s}_{jk} \left( Tr \tilde{\Omega}_j J \check{\Xi}_k^{[1]} J + Tr \check{\Omega}_k^{[1]} J \tilde{\Xi}_j J \right) \right\} \end{aligned} \quad (3.13)$$

Now, we use the superbosonization formulae, i.e., (3.3) and (3.4), to change to the reduced variables as

$$\begin{aligned} \tilde{X}_j &\rightarrow X_j, \quad \tilde{Y}_j \rightarrow Y_j, \quad \tilde{\Omega}_j \rightarrow \Omega_j, \quad \tilde{\Xi}_j \rightarrow \Xi_j, \quad j = 1, \dots, W, \\ \check{X}_k^{[1]} &\rightarrow X_k^{[1]}, \quad \check{\Omega}_k^{[1]} \rightarrow \omega_k^{[1]}(\mathbf{w}_k^{[1]})^*, \quad \check{\Xi}_k^{[1]} \rightarrow \mathbf{w}_k^{[1]}\xi_k^{[1]}, \quad \check{Y}_k^{[1]} \rightarrow Y_k^{[1]} := y_k^{[1]}\mathbf{w}_k^{[1]}(\mathbf{w}_k^{[1]})^*, \quad k = p, q. \end{aligned} \quad (3.14)$$

Here, for  $j = 1, \dots, W$ ,  $X_j$  is a  $2 \times 2$  unitary matrix;  $Y_j$  is a  $2 \times 2$  positive-definite matrix;  $\Omega_j = (\omega_{j,\alpha\beta})$  and  $\Xi_j = (\xi_{j,\alpha\beta})$  are  $2 \times 2$  Grassmann matrices. For  $k = p$  or  $q$ ,  $X_k^{[1]}$  is a  $2 \times 2$  unitary matrix;  $y_k^{[1]}$  is a positive variable;  $\omega_k^{[1]} = (\omega_{k,1}^{[1]}, \omega_{k,2}^{[1]})'$  is a column vector with Grassmann components;  $\xi_k^{[1]} = (\xi_{k,1}^{[1]}, \xi_{k,2}^{[1]})$  is a row vector with Grassmann components. In addition, for  $k = p, q$ ,

$$\mathbf{w}_k^{[1]} = \left( v_k^{[1]}, u_k^{[1]} e^{\mathbf{i}\sigma_k^{[1]}} \right)', \quad u_k^{[1]} = \sqrt{1 - (v_k^{[1]})^2}, \quad v_k^{[1]} \in \mathbb{I}, \quad \sigma_k^{[1]} \in \mathbb{L}. \quad (3.15)$$

Then by using superbosonization formulae, we arrive at the representation

$$\mathbb{E}|G_{pq,11}(z)|^{2n}$$

$$\begin{aligned}
&= \frac{(-1)^W}{(n!)^2 \pi^{2W+4}} \int dX^{[1]} d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\boldsymbol{\omega}^{[1]} d\boldsymbol{\xi}^{[1]} dX dY d\Omega d\Xi \left( y_p^{[1]} y_q^{[1]} (\mathbf{w}_q^{[1]} (\mathbf{w}_q^{[1]})^*)_{12} (\mathbf{w}_p^{[1]} (\mathbf{w}_p^{[1]})^*)_{21} \right)^n \\
&\times \exp \left\{ \mathbf{i} \sum_{j=1}^W \left( \text{Tr} X_j J Z + \text{Tr} Y_j J Z \right) \right\} \exp \left\{ \frac{1}{2M} \sum_{j,k} \tilde{s}_{jk} \left( \text{Tr} X_j J X_k J - \text{Tr} Y_j J Y_k J \right) \right\} \\
&\times \exp \left\{ -\frac{1}{M} \sum_{j,k} \tilde{s}_{jk} \text{Tr} \Omega_j J \Xi_k J \right\} \prod_j \frac{\det^M Y_j}{\det^M (X_j - \Omega_j Y_j^{-1} \Xi_j)} \prod_{k=p,q} \frac{\det (X_k - \Omega_k Y_k^{-1} \Xi_k)}{\det Y_k} \\
&\times \prod_{k=p,q} \exp \left\{ \mathbf{i} \text{Tr} X_k^{[1]} J Z + \mathbf{i} \text{Tr} Y_k^{[1]} J Z \right\} \prod_{k=p,q} \exp \left\{ \frac{1}{M} \sum_{j=1}^W \tilde{s}_{jk} \left( \text{Tr} X_j J X_k^{[1]} J - \text{Tr} Y_j J Y_k^{[1]} J \right) \right\} \\
&\times \prod_{k,\ell=p,q} \exp \left\{ \frac{\tilde{s}_{k\ell}}{2M} \left( \text{Tr} X_k^{[1]} J X_\ell^{[1]} J - \text{Tr} Y_k^{[1]} J Y_\ell^{[1]} J \right) \right\} \prod_{k,\ell=p,q} \exp \left\{ -\frac{\tilde{s}_{k\ell}}{M} \text{Tr} \omega_k^{[1]} (\mathbf{w}_k^{[1]})^* J \mathbf{w}_\ell^{[1]} \boldsymbol{\xi}_\ell^{[1]} J \right\} \\
&\times \prod_{k=p,q} \exp \left\{ -\frac{1}{M} \sum_{j=1}^W \tilde{s}_{jk} \text{Tr} \Omega_j J \mathbf{w}_k^{[1]} \boldsymbol{\xi}_k^{[1]} J \right\} \prod_{k=p,q} \exp \left\{ -\frac{1}{M} \sum_{j=1}^W \tilde{s}_{jk} \text{Tr} \omega_k^{[1]} (\mathbf{w}_k^{[1]})^* J \Xi_j J \right\} \\
&\times \prod_{k=p,q} \frac{y_k^{[1]} \left( y_k^{[1]} - \boldsymbol{\xi}_k^{[1]} (X_k^{[1]})^{-1} \omega_k^{[1]} \right)^2}{\det^2 (X_k^{[1]})}, \tag{3.16}
\end{aligned}$$

where we used the notation  $\mathbf{y}^{[1]} := (y_p^{[1]}, y_q^{[1]})$ ,  $\mathbf{w}^{[1]} := (\mathbf{w}_p^{[1]}, \mathbf{w}_q^{[1]})$ . The differentials in (3.16) are defined by

$$\begin{aligned}
dX^{[1]} &:= \prod_{j=p,q} d\hat{\mu}(X_p^{[1]}) d\hat{\mu}(X_q^{[1]}), & d\mathbf{y}^{[1]} &:= \prod_{j=p,q} \mathbf{1}(y_j^{[1]} > 0) dy_j^{[1]}, \\
d\mathbf{w}^{[1]} &:= \prod_{j=p,q} d\mathbf{w}_j^{[1]} = \prod_{j=p,q} v_j^{[1]} dv_j^{[1]} d\sigma_j^{[1]}, & d\boldsymbol{\omega}^{[1]} d\boldsymbol{\xi}^{[1]} &:= \prod_{\alpha=1,2} \prod_{j=p,q} \omega_{j,\alpha}^{[1]} \xi_{j,\alpha}^{[1]}. \\
dX &:= \prod_{j=1}^W d\hat{\mu}(X_j), & dY &:= \prod_{j=1}^W d\hat{\nu}(Y_j), & d\Omega d\Xi &:= \prod_{\alpha,\beta=1,2} \prod_{j=1}^W d\omega_{j,\alpha\beta} d\xi_{j,\alpha\beta}.
\end{aligned}$$

The regions of the integral of  $X_j$ 's and  $X_k^{[1]}$ 's are all  $U(2)$ , and those of  $Y_j$ 's are the set of all positive-definite matrices. The integral of  $v_k^{[1]}$  ranges over  $\mathbb{I}$  and that of  $\sigma_k^{[1]}$  ranges over  $\mathbb{L}$ , for  $k = p, q$ .

Now we change the variables as  $X_j J \rightarrow X_j$ ,  $Y_j J \rightarrow B_j$ ,  $\Omega_j J \rightarrow \Omega_j$ ,  $\Xi_j J \rightarrow \Xi_j$  and perform the scaling  $X_j \rightarrow -M X_j$ ,  $B_j \rightarrow M B_j$ ,  $\Omega_j \rightarrow \sqrt{M} \Omega_j$  and  $\Xi_j \rightarrow \sqrt{M} \Xi_j$ . Consequently, we can write

$$\begin{aligned}
\mathbb{E}|G_{pq,11}(z)|^{2n} &= \frac{(-1)^W M^{4W}}{(n!)^2 \pi^{2W+4}} \int dX^{[1]} d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\boldsymbol{\omega}^{[1]} d\boldsymbol{\xi}^{[1]} dX dB d\Omega d\Xi \exp \left\{ -M(K(X) + L(B)) \right\} \\
&\times \mathcal{P}(\Omega, \Xi, X, B) \cdot \mathcal{Q}(\Omega, \Xi, \boldsymbol{\omega}^{[1]}, \boldsymbol{\xi}^{[1]}, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \cdot \mathcal{F}(X, B, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}), \tag{3.17}
\end{aligned}$$

where the functions in the integrand are defined as

$$\begin{aligned}
K(X) &:= -\frac{1}{2} \sum_{j,k} \tilde{s}_{jk} \text{Tr} X_j X_k + \mathbf{i} E \sum_j \text{Tr} X_j + \sum_j \log \det X_j, \\
L(B) &:= \frac{1}{2} \sum_{j,k} \tilde{s}_{jk} \text{Tr} B_j B_k - \mathbf{i} E \sum_j \text{Tr} B_j - \sum_j \log \det B_j, \\
\mathcal{P}(\Omega, \Xi, X, B) &:= \exp \left\{ -\sum_{j,k} \tilde{s}_{jk} \text{Tr} \Omega_j \Xi_k \right\} \prod_j \frac{1}{\det^M (1 + M^{-1} X_j^{-1} \Omega_j B_j^{-1} \Xi_j)} \\
&\times \prod_{k=p,q} \frac{\det (X_k + M^{-1} \Omega_k B_k^{-1} \Xi_k)}{\det B_k},
\end{aligned}$$

$$\begin{aligned}
\mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) &:= \prod_{k=p,q} \left(1 - (y_k^{[1]})^{-1} \xi_k^{[1]} (X_k^{[1]})^{-1} \omega_k^{[1]}\right)^2 \\
&\times \prod_{k=p,q} \exp \left\{ -\frac{1}{\sqrt{M}} \sum_j \tilde{\mathbf{s}}_{jk} \left( Tr \Omega_j \mathbf{w}_k^{[1]} \xi_k^{[1]} J + Tr \omega_k^{[1]} (\mathbf{w}_k^{[1]})^* J \Xi_j \right) \right\} \\
&\times \prod_{k,\ell=p,q} \exp \left\{ -\frac{1}{M} \tilde{\mathbf{s}}_{k\ell} Tr \omega_k^{[1]} (\mathbf{w}_k^{[1]})^* J \mathbf{w}_\ell^{[1]} \xi_\ell^{[1]} J \right\}, \\
\mathcal{F}(X, B, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) &:= f(X, X^{[1]}) g(B, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}), \tag{3.18}
\end{aligned}$$

with

$$\begin{aligned}
f(X, X^{[1]}) &:= \exp \left\{ M\eta \sum_{j=1}^W Tr X_j J \right\} \prod_{k,\ell=p,q} \exp \left\{ \frac{\tilde{\mathbf{s}}_{k\ell}}{2M} Tr X_k^{[1]} J X_\ell^{[1]} J \right\} \\
&\times \prod_{k=p,q} \frac{1}{\det^2(X_k^{[1]})} \exp \left\{ i Tr X_k^{[1]} J Z - \sum_j \tilde{\mathbf{s}}_{jk} Tr X_j X_k^{[1]} J \right\}, \tag{3.19}
\end{aligned}$$

$$\begin{aligned}
g(B, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) &:= \exp \left\{ -M\eta \sum_{j=1}^W Tr B_j J \right\} \prod_{k,\ell=p,q} \exp \left\{ -\frac{\tilde{\mathbf{s}}_{k\ell}}{2M} Tr Y_k^{[1]} J Y_\ell^{[1]} J \right\} \\
&\times \left( (\mathbf{w}_q^{[1]} (\mathbf{w}_q^{[1]})^*)_{12} (\mathbf{w}_p^{[1]} (\mathbf{w}_p^{[1]})^*)_{21} \right)^n \prod_{k=p,q} (y_k^{[1]})^{n+3} \exp \left\{ i Tr Y_k^{[1]} J Z - \sum_j \tilde{\mathbf{s}}_{jk} Tr B_j Y_k^{[1]} J \right\}. \tag{3.20}
\end{aligned}$$

In (3.17), the regions of  $X_j$ 's and  $X_k^{[1]}$ 's are all  $U(2)$ , and those of  $B_j$ 's are the set of the matrices  $A$  satisfying  $AJ > 0$ . Roughly speaking, here we collected the terms containing the Grassmann variables from  $\omega_k^{[1]}$ 's and  $\xi_k^{[1]}$ 's, resulting the factor  $\mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]})$ . Then, we put the terms containing the Grassmann variables in  $\Omega_j$ 's and  $\Xi_j$ 's together (except those in  $\mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]})$ ), resulting the factor  $\mathcal{P}(\Omega, \Xi, X, B)$ . Finally, we sorted out the terms containing  $\eta$ ,  $y_k^{[1]}$  and the variables from  $X_k^{[1]}$ 's and  $\mathbf{w}_k^{[1]}$ 's (except those in  $\mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]})$ ), resulting the factor  $\mathcal{F}(X, B, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]})$ . This separation indicates the order in which we will perform the integrations.

**3.6. Parametrization for  $X, B$ .** Similarly to the discussion in [19], we start with some preliminary parameterization. At first, we do the eigendecomposition

$$X_j = P_j^* \hat{X}_j P_j, \quad B_j = Q_j^{-1} \hat{B}_j Q_j, \quad P_j \in \dot{U}(2), \quad Q_j \in \dot{U}(1, 1), \tag{3.21}$$

where

$$\hat{X}_j = \text{diag}(x_{j,1}, x_{j,2}), \quad \hat{B}_j = \text{diag}(b_{j,1}, -b_{j,2}), \quad x_{j,1}, x_{j,2} \in \Sigma, \quad b_{j,1}, b_{j,2} \in \mathbb{R}_+. \tag{3.22}$$

Further, we introduce

$$V_j = P_j P_1^* \in \dot{U}(2), \quad T_j = Q_j Q_1^{-1} \in \dot{U}(1, 1), \quad j = 1, \dots, W. \tag{3.23}$$

Especially, we have  $V_1 = T_1 = I$ . Now, we parameterize  $P_1, Q_1, V_j$  and  $T_j$  for all  $j = 2, \dots, W$  as follows

$$\begin{aligned}
P_1 &= \begin{pmatrix} u & v e^{i\theta} \\ -v e^{-i\theta} & u \end{pmatrix}, \quad V_j = \begin{pmatrix} u_j & v_j e^{i\theta_j} \\ -v_j e^{-i\theta_j} & u_j \end{pmatrix}, \\
u &= \sqrt{1-v^2}, \quad u_j = \sqrt{1-v_j^2}, \quad v, v_j \in \mathbb{I}, \quad \theta, \theta_j \in \mathbb{L}, \\
Q_1 &= \begin{pmatrix} s & t e^{i\sigma} \\ t e^{-i\sigma} & s \end{pmatrix}, \quad T_j = \begin{pmatrix} s_j & t_j e^{i\sigma_j} \\ t_j e^{-i\sigma_j} & s_j \end{pmatrix},
\end{aligned}$$

$$s = \sqrt{1+t^2}, \quad s_j = \sqrt{1+t_j^2}, \quad t, t_j \in \mathbb{R}_+, \quad \sigma, \sigma_j \in \mathbb{L}. \quad (3.24)$$

Under the parametrization above, we can express the corresponding differentials as follows.

$$\begin{aligned} dX dB &= d\mu(P_1) d\nu(Q_1) \cdot \prod_{j=2}^W d\mu(V_j) d\nu(T_j) \cdot \prod_{j=1}^W db_{j,1} db_{j,2} \cdot \frac{dx_{j,1}}{2\pi i} \frac{dx_{j,2}}{2\pi i} \\ &\quad \times 2^W (\pi/2)^{2W} \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2, \end{aligned} \quad (3.25)$$

where

$$d\mu(P_1) = 2v dv \cdot \frac{d\theta}{2\pi}, \quad d\mu(V_j) = 2v_j dv_j \cdot \frac{d\theta_j}{2\pi}, \quad d\nu(Q_1) = 2t dt \cdot \frac{d\sigma}{2\pi}, \quad d\nu(T_j) = 2t_j dt_j \cdot \frac{d\sigma_j}{2\pi}.$$

In addition, for simplicity, we do the change of variables

$$\Omega_j \rightarrow P_1^* \Omega_j Q_1, \quad \Xi_j \rightarrow Q_1^{-1} \Xi_j P_1. \quad (3.26)$$

Note that the Berezinian of such a change is 1. After this change,  $\mathcal{P}(\Omega, \Xi, X, B, \mathbf{y}^{[1]}, \mathbf{w}^{[1]})$  turns out to be independent of  $P_1$  and  $Q_1$ .

To adapt to the new parametrization, we change the notation

$$\begin{aligned} K(X) &\rightarrow K(\hat{X}, V), \quad L(B) \rightarrow L(\hat{B}, T), \quad \mathcal{P}(\Omega, \Xi, X, B) \rightarrow \mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T), \\ \mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) &\rightarrow \mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}), \\ \mathcal{F}(X, B, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) &\rightarrow \mathcal{F}(\hat{X}, \hat{B}, V, T, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}), \\ f(X, X^{[1]}) &\rightarrow f(P_1, V, \hat{X}, X^{[1]}), \quad g(B, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \rightarrow g(Q_1, T, \hat{B}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}). \end{aligned} \quad (3.27)$$

We recall here that  $K(X)$  does not depend on  $P_1$ , as well,  $L(B)$  does not depend on  $Q_1$ . Moreover, according to the change (3.26), we have

$$\begin{aligned} \mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T) &= \exp \left\{ - \sum_{j,k} \tilde{s}_{jk} Tr \Omega_j \Xi_k \right\} \cdot \prod_j \frac{1}{\det^M (1 + M^{-1} V_j^* \hat{X}_j^{-1} V_j \Omega_j T_j^{-1} \hat{B}_j^{-1} T_j \Xi_j)} \\ &\quad \times \prod_{k=p,q} \frac{\det (V_k^* \hat{X}_k V_k + M^{-1} \Omega_k T_k^{-1} \hat{B}_k^{-1} T_k \Xi_k)}{\det \hat{B}_k} \end{aligned} \quad (3.28)$$

and

$$\begin{aligned} &\mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \\ &= \prod_{k=p,q} \exp \left\{ - \frac{1}{\sqrt{M}} \sum_j \tilde{s}_{jk} \left( Tr P_1^* \Omega_j Q_1 \mathbf{w}_k^{[1]} \xi_k^{[1]} J + Tr \omega_k^{[1]} (\mathbf{w}_k^{[1]})^* J Q_1^{-1} \Xi_j P_1 \right) \right\} \\ &\quad \times \prod_{k,\ell=p,q} \exp \left\{ - \frac{1}{M} \tilde{s}_{k\ell} Tr \omega_k^{[1]} (\mathbf{w}_k^{[1]})^* J (\mathbf{w}_\ell^{[1]}) \xi_\ell^{[1]} J \right\} \cdot \prod_{k=p,q} \left( 1 - (y_k^{[1]})^{-1} \xi_k^{[1]} (X_k^{[1]})^{-1} \omega_k^{[1]} \right)^2. \end{aligned} \quad (3.29)$$

Consequently, using (3.25), from (3.17) we can write

$$\begin{aligned} \mathbb{E} |G_{pq,11}(z)|^{2n} &= \frac{M^{4W}}{(n!)^2 8^W \pi^{2W+4}} \int \prod_{j=2}^W d\mu(V_j) d\nu(T_j) \int_{\mathbb{R}_+^{2W}} \prod_{j=1}^W db_{j,1} db_{j,2} \oint_{\Sigma^{2W}} \prod_{j=1}^W dx_{j,1} dx_{j,2} \\ &\quad \times \exp \left\{ -M(K(\hat{X}, V) + L(\hat{B}, T)) \right\} \cdot \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot A(\hat{X}, \hat{B}, V, T). \end{aligned} \quad (3.30)$$

where we introduced the notation

$$\begin{aligned} \mathbf{A}(\hat{X}, \hat{B}, V, T) &:= \int dX^{[1]} d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\boldsymbol{\omega}^{[1]} d\boldsymbol{\xi}^{[1]} d\Omega d\Xi d\mu(P_1) d\nu(Q_1) \mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T) \\ &\times \mathcal{Q}(\Omega, \Xi, \boldsymbol{\omega}^{[1]}, \boldsymbol{\xi}^{[1]}, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \cdot \mathcal{F}(\hat{X}, \hat{B}, V, T, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}). \end{aligned} \quad (3.31)$$

In (3.30), the regions of  $V_j$ 's are all  $\mathring{U}(2)$ , and those of  $T_j$ 's are all  $\mathring{U}(1, 1)$ . Observe that all Grassmann variables are inside the integrand of the integral  $\mathbf{A}(\hat{X}, \hat{B}, V, T)$ . Hence, (3.30) separates the saddle point calculation from the observable  $\mathbf{A}(\hat{X}, \hat{B}, V, T)$ .

To facilitate the discussions in the remaining part, we introduce some additional terms and notation here. Henceforth, we will employ the notation

$$(X^{[1]})^{-1} = \left\{ (X_p^{[1]})^{-1}, (X_q^{[1]})^{-1} \right\}, \quad (\mathbf{y}^{[1]})^{-1} = \left\{ (y_p^{[1]})^{-1}, (y_q^{[1]})^{-1} \right\}$$

for the collection of inverse matrices and reciprocals, respectively. For a matrix or a vector  $A$  under discussion, we will use the term *A-variables* to refer to all the variables parametrizing it. For example,  $\hat{X}_j$ -variables means  $x_{j,1}$  and  $x_{j,2}$ , and  $\hat{X}$ -variables refer to the collection of all  $\hat{X}_j$ -variables. Analogously, we can define the terms *T-variables*,  $\mathbf{y}^{[1]}$ -variables,  $\Omega$ -variables and so on. We use another term *A-entries* to refer to the non-zero entries of  $A$ . Note that  $\hat{X}_j$ -variables are just  $\hat{X}_j$ -entries. However, for  $T_j$ , they are different, namely,

$$T_j\text{-variables: } t_j, \sigma_j, \quad \text{vs. } T_j\text{-entries: } s_j, t_j e^{i\sigma_j}, t_j e^{-i\sigma_j}.$$

Analogously, we will also use the term *T-entries* to refer to the collection of all  $T_j$ -entries. Then  $V$ -entries,  $\mathbf{w}^{[1]}$ -entries, etc. are defined in the same manner. It is easy to check that  $Q_1^{-1}$ -entries are the same as  $Q_1$ -entries, up to a sign, as well,  $T_j^{-1}$ -entries are the same as  $T_j$ -entries, for all  $j = 2, \dots, W$ .

Moreover, to simplify the notation, we make the convention here that we will frequently use a *dot* to represent all the arguments of a function. That means, for instance, we will write  $\mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T)$  as  $\mathcal{P}(\cdot)$  if there is no confusion. Analogously, we will also use the abbreviation  $\mathcal{Q}(\cdot)$ ,  $\mathcal{F}(\cdot)$ ,  $\mathbf{A}(\cdot)$ , and so on.

Let  $\mathbf{a} := \{a_1, \dots, a_\ell\}$  be a set of variables, we will adopt the notation

$$\mathfrak{Q}(\mathbf{a}; \kappa_1, \kappa_2, \kappa_3)$$

to denote the class of all multivariate polynomials  $\mathbf{p}(\mathbf{a})$  in the arguments  $a_1, \dots, a_\ell$  such that the following three conditions are satisfied: (i) The total number of the monomials in  $\mathbf{p}(\mathbf{a})$  is bounded by  $\kappa_1$ ; (ii) the coefficients of all monomials in  $\mathbf{p}(\mathbf{a})$  are bounded by  $\kappa_2$  in magnitude; (iii) the power of each  $a_i$  in each monomial is bounded by  $\kappa_3$ , for all  $i = 1, \dots, \ell$ . For example,

$$5b_{j,1}^{-1} + 3b_{j,1}t_j^2 + 1 \in \mathfrak{Q}(\{b_{j,1}^{-1}, b_{j,1}, t_j\}; 3, 5, 2). \quad (3.32)$$

In addition, we define the subset of  $\mathfrak{Q}(\mathbf{a}; \kappa_1, \kappa_2, \kappa_3)$ , namely,

$$\mathfrak{Q}_{\deg}(\mathbf{a}; \kappa_1, \kappa_2, \kappa_3)$$

consisting of those polynomials in  $\mathfrak{Q}(\mathbf{a}; \kappa_1, \kappa_2, \kappa_3)$  such that the degree is bounded by  $\kappa_3$ , i.e. the total degree of each monomial is bounded by  $\kappa_3$ . For example

$$5b_{j,1}^{-1} + 3b_{j,1}t_j^2 + 1 \in \mathfrak{Q}_{\deg}(\{b_{j,1}^{-1}, b_{j,1}, t_j\}; 3, 5, 3).$$

## 4. PRELIMINARY DISCUSSION ON THE INTEGRAND

In this section, we perform a preliminary analysis on the factors of the integrand in (3.17). For convenience, we introduce the matrix

$$\mathfrak{J} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}. \quad (4.1)$$

4.1.  $\exp\{-M(K(\hat{X}, V) + L(\hat{B}, T))\}$ . Recall the parametrization of  $\hat{B}_j$ ,  $\hat{X}_j$ ,  $T_j$  and  $V_j$  in (3.22) and (3.24), as well as the matrices defined in (1.24). According to the discussion in [19], there are three types of saddle points of this function, namely,

- Type I : For each  $j$ ,  $(\hat{B}_j, T_j, \hat{X}_j) = (D_{\pm}, I, D_{\pm})$  or  $(D_{\pm}, I, D_{\mp})$ ,  
 $\theta_j \in \mathbb{L}$ ,  $v_j = 0$  if  $\hat{X}_j = \hat{X}_1$ , and  $v_j = 1$  if  $\hat{X}_j \neq \hat{X}_1$ .
- Type II : For each  $j$ ,  $(\hat{B}_j, T_j, \hat{X}_j) = (D_{\pm}, I, D_{+})$  and  $V_j \in \mathring{U}(2)$ .
- Type III : For each  $j$ ,  $(\hat{B}_j, T_j, \hat{X}_j) = (D_{\pm}, I, D_{-})$  and  $V_j \in \mathring{U}(2)$ .

(Actually, since  $\theta_j$  and  $v_j$  vary on continuous sets, it would be more appropriate to use the term *saddle manifolds*.) Note that at each type of saddle points, we have  $(\hat{B}_j, T_j) = (D_{\pm}, I)$  for all  $j$ . We will see that the main contribution to the integral (3.17) comes from some small vicinities of the Type I saddle points. Furthermore, the contributions from all the Type I saddle points are the same, which can be explained as follows. At first, by the definition in (3.23), we have  $V_1 = I$ . If we regard  $\theta_j$ 's in the parametrization of  $V_j$ 's as fixed parameters, it is easy to see that there are totally  $2^W$  choices of Type I saddle points. Moreover, if  $v_j = 1$ , we can do the transform

$$\hat{X}_j \rightarrow \mathfrak{J}\hat{X}_j\mathfrak{J} = \hat{X}_1, \quad V_j \rightarrow \mathfrak{J}V_j = I \quad \text{in } \mathring{U}(1, 1).$$

Consequently, it suffices to consider two saddle points

$$(\hat{B}_j, T_j, \hat{X}_j, V_j) = (D_{\pm}, I, D_{\pm}, I), \quad \text{or} \quad (D_{\pm}, I, D_{\mp}, I), \quad (4.2)$$

corresponding to  $\hat{X}_1 = D_{\pm}$  or  $D_{\mp}$ , respectively. Furthermore, the contributions to the integral (3.17) from the vicinities of these two saddle points are also the same. To see this, we recall the fact that the original integrand in (3.17) is a function of the entries of  $X_j = P_j^{-1}\hat{X}_jP_j$ . Now we do the transform  $P_j \rightarrow \mathfrak{J}P_j$  and  $\hat{X}_j \rightarrow \mathfrak{J}\hat{X}_j\mathfrak{J}$  for all  $j = 1, \dots, W$  to change one saddle in (4.2) to the other. Now, since the Haar measure on  $\mathring{U}(2)$  is invariant under the shift  $P_1 \rightarrow \mathfrak{J}P_1$ , the integral over  $P_1$ -variables is unchanged. That means, for Type I saddle points, it suffices to consider

- Type I' : For each  $j$ ,  $(\hat{B}_j, T_j, \hat{X}_j, V_j) = (D_{\pm}, I, D_{\pm}, I)$ .

In summary, the total contribution to the integral (3.17) from all Type I saddle points is  $2^W$  times that from the Type I' saddle point.

Following the discussion in [19], we will show in Section 5 that both  $K(\hat{X}, V) - K(D_{\pm}, I)$  and  $L(\hat{B}, T) - L(D_{\pm}, I)$  have positive real parts, bounded by some positive quadratic forms from below, which allows us to perform the saddle point analysis. In addition, it will be seen that in a vicinity of Type I' saddle point,  $\exp\{-M(K(\hat{X}, V) + L(\hat{B}, T))\}$  is approximately Gaussian.

4.2.  $\mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]})$ . The function  $\mathcal{Q}(\cdot)$  contains both the  $\Omega, \Xi$ -variables from  $\mathcal{P}(\cdot)$ , and the  $P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}$ -variables from  $\mathcal{F}(\cdot)$ . In addition, note that in the integrand in (3.17),  $\mathcal{Q}(\cdot)$  is the only factor containing the  $\omega^{[1]}$  and  $\xi^{[1]}$ -variables. Hence, we can compute the integral

$$\mathcal{Q}(\Omega, \Xi, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) := \int d\omega^{[1]} d\xi^{[1]} \mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \quad (4.3)$$

at first. The explicit formula for  $\mathcal{Q}(\cdot)$  is complicated and irrelevant for us. From (3.29) and the definition of the Grassmann integral, it is not difficult to see that  $\mathcal{Q}(\cdot)$  is a polynomial of the  $(X^{[1]})^{-1}$ ,  $(\mathbf{y}^{[1]})^{-1}$ ,

$\mathbf{w}^{[1]}$ ,  $P_1$ ,  $Q_1$ ,  $\Omega$  and  $\Xi$ -entries. In principle, for each monomial in the polynomial  $Q(\cdot)$ , we can combine the Grassmann variables with  $\mathcal{P}(\cdot)$ , then perform the integral over  $\Omega$  and  $\Xi$ -variables, whilst we combine the complex variables with  $\mathcal{F}(\cdot)$ , and perform the integral over  $X^{[1]}$ ,  $\mathbf{y}^{[1]}$ ,  $\mathbf{w}^{[1]}$ ,  $P_1$  and  $Q_1$ -variables. A formal discussion on  $Q(\cdot)$  will be given in Section 6.1. However, the terms from  $Q(\cdot)$  turn out to be irrelevant in our proof. Therefore, in the arguments with  $Q(\cdot)$  involved, a typical strategy that we will adopt is as follows: we usually neglect  $Q(\cdot)$  at first, and perform the discussion on  $\mathcal{P}(\cdot)$  and  $\mathcal{F}(\cdot)$  separately, at the end, we make necessary comments on how to slightly modify the discussions to take  $Q(\cdot)$  into account.

4.3.  $\mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T)$ . We will mainly regard  $\mathcal{P}(\cdot)$  as a function of the  $\Omega$  and  $\Xi$ -variables. As mentioned above, we also have some  $\Omega$  and  $\Xi$ -variables from the irrelevant term  $Q(\cdot)$ . But we temporarily ignore them and regard as if the integral over  $\Omega$  and  $\Xi$ -variables reads

$$P(\hat{X}, \hat{B}, V, T) := \int d\Omega d\Xi \mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T). \quad (4.4)$$

We shall estimate  $P(\cdot)$  in three different regions: (1) the complement of the vicinities of the saddle points; (2) the vicinity of Type I saddle point; (3) the vicinities of Type II and III saddle points, which will be done in Sections 6.2, 9.1 and 10.1, respectively. (Definition 5.5 gives the precise definition of the vicinities.) In each case we will decompose the function  $\mathcal{P}(\cdot)$  as a product of a Gaussian measure and a multivariate polynomial of Grassmann variables. Consequently, we can employ (3.2) to perform the integral of this polynomial against the Gaussian measure, whereby  $P(\cdot)$  can be estimated. Especially, it turns out that in the vicinity of Type I saddle points,  $P(\cdot)$  is approximately the normalizing constant of the Gaussian measure obtained from  $\exp\{-M(K(\hat{X}, V) + L(\hat{B}, T))\}$ .

4.4.  $\mathcal{F}(\hat{X}, \hat{B}, V, T, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]})$ . Observe that  $\mathcal{F}$  is the only term containing the energy scale  $\eta$ . As in the previous discussion of  $\mathcal{P}(\cdot)$ , here we also ignore the  $P_1$ ,  $Q_1$ ,  $X^{[1]}$ ,  $\mathbf{y}^{[1]}$ ,  $\mathbf{w}^{[1]}$ -variables from the irrelevant term  $Q(\cdot)$  temporarily, and investigate the integral

$$\begin{aligned} F(\hat{X}, \hat{B}, V, T) &= \int dX^{[1]} d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\mu(P_1) d\nu(Q_1) \mathcal{F}(\hat{X}, \hat{B}, V, T, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \\ &= \int dX^{[1]} d\mu(P_1) f(\hat{X}, V, P_1) \cdot \int d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\nu(Q_1) g(\hat{B}, T, Q_1, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}). \end{aligned} \quad (4.5)$$

We shall also estimate  $F(\cdot)$  in three different regions: (1) the complement of the vicinities of the saddle points; (2) the vicinity of Type I saddle point; (3) the vicinities of Type II and III saddle points, which will be done in Sections 6.3, 9.2 and 10.2, respectively.

Especially, when we restrict the  $\hat{X}$ ,  $\hat{B}$ ,  $V$  and  $T$ -variables to the vicinity of the Type I saddle points, the above integral can be performed approximately, resulting our main term, a factor of order  $1/(N\eta)^{n+2}$ . This step will be done in Section 9. It is instructive to give a heuristic sketch of this calculation. At first, we plug the Type I saddle points into (4.5). We will show that the integral of  $f(\cdot)$  approximately reads

$$e^{-(a_+ - a_-)N\eta} \int dX^{[1]} d\mu(P_1) f(D_{\pm}, I, P_1) \sim \frac{1}{N\eta},$$

which is the easy part. Then, recalling the definition of  $g(\cdot)$  in (3.20) and the parameterization (3.15), we will show that the integral of  $g(\cdot)$  approximately reads

$$\begin{aligned} &e^{(a_+ - a_-)N\eta} \int d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\nu(Q_1) g(D_{\pm}, I, Q_1, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \\ &\sim \int_0^\infty 2tdt \int_{\mathbb{L}^2} d\sigma_p^{[1]} d\sigma_q^{[1]} e^{i\eta\sigma_p^{[1]}} e^{-i\eta\sigma_q^{[1]}} \cdot e^{-cN\eta t^2 + c_1 e^{-i\sigma_p^{[1]}} t + c_2 e^{i\sigma_q^{[1]}} t} \\ &\sim \int_0^\infty 2tdt \cdot t^{2n} \cdot e^{-cN\eta t^2} \sim \frac{1}{(N\eta)^{n+1}}, \end{aligned} \quad (4.6)$$



where in the second step above we used the fact that

$$\int_{\mathbb{L}} d\sigma \cdot e^{i\sigma} e^{ce^{-i\sigma}t} \sim t^n.$$

We notice that the factor  $e^{in\sigma_p^{[1]}} e^{-in\sigma_q^{[1]}}$  in (4.6) actually comes from the term

$$\left( (\mathbf{w}_q^{[1]}(\mathbf{w}_q^{[1]})^*)_{12} (\mathbf{w}_p^{[1]}(\mathbf{w}_p^{[1]})^*)_{21} \right)^n$$

in (3.20). This factor brings a strong oscillation to the integrand in the integral (4.6). In Case 2, an analogous factor will appear, resulting the same estimate as (4.6). However, in Case 3, such an oscillating factor is absent, then the estimate for the counterpart of the integral in (4.6) is of order  $1/N\eta$  instead of  $1/(N\eta)^{n+1}$ . The detailed analysis will be presented in Sections 10 and 11.

## 5. SADDLE POINTS AND VICINITIES

In this section, we study the saddle points of  $K(\hat{X}, V)$  and  $L(\hat{B}, T)$  and deform the contours of the  $\hat{B}$ -variables to pass through the saddle points. Then we introduce and classify some small vicinities of these saddle points. The derivation of the saddle points of  $K(\hat{X}, V)$  and  $L(\hat{B}, T)$  in Section 5.1 and 5.2 below is essentially the same as the counterpart in [19], the only difference is that we are working under a more general setting on  $S$ . Hence, in Section 5.1 and 5.2, we just sketch the discussion, list the results, and make necessary modifications to adapt to our setting. In the sequel, we employ the notation

$$\begin{aligned} \mathbf{b}_a &:= (b_{1,a}, \dots, b_{W,a}), \quad \mathbf{x}_a := (x_{1,a}, \dots, x_{W,a}), \quad a = 1, 2, \\ \mathbf{t} &:= (t_2, \dots, t_W), \quad \mathbf{v} := (v_2, \dots, v_W), \quad \boldsymbol{\sigma} := (\sigma_2, \dots, \sigma_W), \quad \boldsymbol{\theta} := (\theta_2, \dots, \theta_W). \end{aligned} \quad (5.1)$$

As mentioned above, later we also need to deform the contours, and discuss the integral over some vicinities of the saddle points, thus it is convenient to introduce a notation for the integral over specific domains. To this end, for  $a = 1, 2$ , we use  $\mathbf{I}_a^b$  and  $\mathbf{I}_a^x$  to denote generic domains of  $\mathbf{b}_a$  and  $\mathbf{x}_a$  respectively. Analogously, we use  $\mathbf{I}^t$  and  $\mathbf{I}^v$  to represent generic domains of  $\mathbf{t}$  and  $\mathbf{v}$ , respectively. These domains will be specified later. Now, for some collection of domains, we introduce the notation

$$\begin{aligned} \mathcal{I}(\mathbf{I}_1^b, \mathbf{I}_2^b, \mathbf{I}_1^x, \mathbf{I}_2^x, \mathbf{I}^t, \mathbf{I}^v) &:= \frac{M^{4W}}{(n!)^{28W} \pi^{2W+4}} \int_{\mathbb{L}^{2W-2}} \prod_{j=2}^W \frac{d\theta_j}{2\pi} \prod_{j=2}^W \frac{d\sigma_j}{2\pi} \\ &\times \int_{\mathbf{I}_1^b} \prod_{j=1}^W db_{j,1} \int_{\mathbf{I}_2^b} \prod_{j=1}^W db_{j,2} \int_{\mathbf{I}_1^x} \prod_{j=1}^W dx_{j,1} \int_{\mathbf{I}_2^x} \prod_{j=1}^W dx_{j,2} \int_{\mathbf{I}^t} \prod_{j=2}^W 2t_j dt_j \int_{\mathbf{I}^v} \prod_{j=2}^W 2v_j dv_j \\ &\times \exp \left\{ -M(K(\hat{X}, V) + L(\hat{B}, T)) \right\} \cdot \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot A(\hat{X}, \hat{B}, V, T). \end{aligned} \quad (5.2)$$

For example, we can write (3.30) as

$$\mathbb{E}|G_{pq,11}(z)|^{2n} = \mathcal{I}(\mathbb{R}_+^W, \mathbb{R}_+^W, \Sigma^W, \Sigma^W, \mathbb{R}_+^{W-1}, \mathbb{I}^{W-1}), \quad (5.3)$$

which is the integral over the full domain.

**5.1. Saddle points of  $L(\hat{B}, T)$ .** We introduce the function

$$\mathbb{k}(a) := \frac{a^2}{2} - \mathbf{i}Ea - \log a, \quad a \in \mathbb{C}. \quad (5.4)$$

Recalling the definition of  $L(\cdot)$  in (3.18), the decomposition of  $B_j$ 's in (3.21) and the definition of  $T_j$ 's in (3.23), we can write

$$L(\hat{B}, T) = -\frac{1}{4} \sum_{j,k} \mathfrak{s}_{jk} \text{Tr} (T_j^{-1} \hat{B}_j T_j - T_k^{-1} \hat{B}_k T_k)^2 + \sum_j \left( \frac{1}{2} \text{Tr} \hat{B}_j^2 - \mathbf{i}E \text{Tr} \hat{B}_j - \log \det \hat{B}_j \right)$$

$$\begin{aligned}
&= -\frac{1}{4} \sum_{j,k} \mathfrak{s}_{jk} \text{Tr}(\hat{B}_j - \hat{B}_k)^2 + \sum_j \left( \frac{1}{2} \text{Tr} \hat{B}_j^2 - \mathbf{i} E \text{Tr} \hat{B}_j - \log \det \hat{B}_j \right) \\
&\quad + \frac{1}{2} \sum_{j,k} \mathfrak{s}_{jk} |(T_k T_j^{-1})_{12}|^2 (b_{j,1} + b_{j,2})(b_{k,1} + b_{k,2}) \\
&=: \ell(\mathbf{b}_1) + \ell(-\mathbf{b}_2) + \ell_S(\hat{B}, T),
\end{aligned} \tag{5.5}$$

where we used the notation introduced in (5.1), and the functions  $\ell(\cdot)$  and  $\ell_S(\cdot)$  are defined as

$$\begin{aligned}
\ell(\mathbf{a}) &:= -\frac{1}{4} \sum_{j,k} \mathfrak{s}_{jk} (a_j - a_k)^2 + \sum_j \mathbb{K}(a_j), \quad \mathbf{a} = (a_1, \dots, a_W) \in \mathbb{C}^W, \\
\ell_S(\hat{B}, T) &:= \frac{1}{2} \sum_{j,k} \mathfrak{s}_{jk} |(T_k T_j^{-1})_{12}|^2 (b_{j,1} + b_{j,2})(b_{k,1} + b_{k,2}).
\end{aligned} \tag{5.6}$$

Following the discussion in [19] with slight modification (see Section 3 therein), we see that for  $|E| \leq \sqrt{2} - \kappa$ , the saddle point of  $L(\hat{B}, T)$  is

$$(\hat{B}_j, T_j) = (D_\pm, I), \quad \forall j = 1, \dots, W, \tag{5.7}$$

where  $D_\pm$  is defined in (1.24). For simplicity, we will write (5.7) as  $(\hat{B}, T) = (D_\pm, I)$  in the sequel. Observe that

$$L(D_\pm, I) = \ell(a_+) + \ell(a_-), \quad \ell_S(D_\pm, I) = 0. \tag{5.8}$$

We introduce the notation

$$\mathring{\ell}_{++}(\mathbf{a}) := \ell(\mathbf{a}) - \ell(a_+), \quad \mathring{\ell}_{+-}(\mathbf{a}) := \ell(\mathbf{a}) - \ell(a_-), \quad \mathring{\ell}_{--}(\mathbf{a}) := \ell(-\mathbf{a}) - \ell(a_-), \tag{5.9}$$

where  $\ell(a_+)$  represents the value of  $\ell(\mathbf{a})$  at the point  $\mathbf{a} = (a_+, \dots, a_+)$ , and  $\ell(a_-)$  is defined analogously. Correspondingly, we adopt the notation

$$\mathring{L}(\hat{B}, T) := L(\hat{B}, T) - L(D_\pm, I) = \mathring{\ell}_{++}(\mathbf{b}_1) + \mathring{\ell}_{--}(\mathbf{b}_2) + \ell_S(\hat{B}, T), \tag{5.10}$$

which is implied by (5.5), (5.8) and (5.9). Now, for each  $j = 1, \dots, W$ , we deform the contours of  $b_{j,1}$  and  $b_{j,2}$  to

$$b_{j,1} \in \Gamma := \{ra_+ | r \in \mathbb{R}_+\}, \quad b_{j,2} \in \bar{\Gamma} = \{-ra_- | r \in \mathbb{R}_+\} \tag{5.11}$$

to pass through the saddle points of  $\hat{B}$ -variables, based on the following lemma which will be proved in Section 7.

**Lemma 5.1.** *With the notation introduced in (5.2), we have*

$$\mathcal{I}(\Gamma^W, \bar{\Gamma}^W, \Sigma^W, \Sigma^W, \mathbb{R}_+^{W-1}, \mathbb{I}^{W-1}) = \mathcal{I}(\mathbb{R}_+^W, \mathbb{R}_+^W, \Sigma^W, \Sigma^W, \mathbb{R}_+^{W-1}, \mathbb{I}^{W-1}) = \mathbb{E}|G_{pq,11}(z)|^{2n}.$$

We introduce the notation

$$r_{j,1} = |b_{j,1}|, \quad r_{j,2} = |b_{j,2}|, \quad j = 1, \dots, W. \tag{5.12}$$

Along the new contours, we have the following lemma.

**Lemma 5.2.** *Suppose that  $|E| \leq \sqrt{2} - \kappa$ . Let  $\mathbf{b}_1 \in \Gamma^W$ ,  $\mathbf{b}_2 \in \bar{\Gamma}^W$ . We have*

$$\text{Re} \mathring{L}(\hat{B}, T) \geq c \sum_{a=1,2} \sum_{j=1}^W \left( (r_{j,a} - 1)^2 + (r_{j,a} - \log r_{j,a} - 1) \right) + \text{Re} \ell_S(\hat{B}, T) \geq c \sum_{a=1,2} \sum_{j=1}^W (r_{j,a} - 1)^2 \tag{5.13}$$

for some positive constant  $c$ .

*Proof.* Since  $|E| \leq \sqrt{2} - \kappa$ , we have  $\operatorname{Re}(b_{j,1} + b_{j,2})(b_{k,1} + b_{k,2}) \geq 0$  for  $\mathbf{b}_1 \in \Gamma^W$  and  $\mathbf{b}_2 \in \bar{\Gamma}^W$ , thus  $\operatorname{Re} \ell_S(\hat{B}, T) \geq 0$ , in light of the definition in (5.6). Consequently, according to (5.10), it suffices to prove

$$\operatorname{Re} \ell_{++}(\mathbf{b}_1) + \operatorname{Re} \ell_{--}(\mathbf{b}_2) \geq c \sum_{a=1,2} \sum_{j=1}^W \left( (r_{j,a} - 1)^2 + (r_{j,a} - \log r_{j,a} - 1) \right) \quad (5.14)$$

for some positive constant  $c$ . To see this, we observe the following identities obtained via elementary calculation,

$$\begin{aligned} \operatorname{Re} \ell_{++}(\mathbf{b}_1) &= \frac{E^2 - 2}{4} \left( \frac{1}{2} \sum_{j,k} \mathfrak{s}_{jk} (r_{j,1} - r_{k,1})^2 - \sum_j (r_{j,1} - 1)^2 \right) + \sum_{j=1}^W (r_{j,1} - \log r_{j,1} - 1) \\ \operatorname{Re} \ell_{--}(\mathbf{b}_2) &= \frac{E^2 - 2}{4} \left( \frac{1}{2} \sum_{j,k} \mathfrak{s}_{jk} (r_{j,2} - r_{k,2})^2 - \sum_j (r_{j,2} - 1)^2 \right) + \sum_{j=1}^W (r_{j,2} - \log r_{j,2} - 1), \end{aligned}$$

which together with  $|E| \leq \sqrt{2} - \kappa$  and (1.4) implies (5.14) immediately. Hence, we completed the proof of Lemma 5.2.  $\square$

**5.2. Saddle points of  $K(\hat{X}, V)$ .** Analogously, recalling the definition in (5.6), we can write

$$\begin{aligned} K(\hat{X}, V) &= \frac{1}{4} \sum_{j,k} \mathfrak{s}_{jk} \operatorname{Tr} (V_j^* \hat{X}_j V_j - V_k^* \hat{X}_k V_k)^2 - \sum_j \left( \frac{1}{2} \operatorname{Tr} \hat{X}_j^2 - \mathbf{i} E \operatorname{Tr} \hat{X}_j - \log \det \hat{X}_j \right) \\ &= \frac{1}{4} \sum_{j,k} \mathfrak{s}_{jk} \operatorname{Tr} (\hat{X}_j - \hat{X}_k)^2 - \sum_j \left( \frac{1}{2} \operatorname{Tr} \hat{X}_j^2 - \mathbf{i} E \operatorname{Tr} \hat{X}_j - \log \det \hat{X}_j \right) \\ &\quad + \frac{1}{2} \sum_{j,k} \mathfrak{s}_{jk} |(V_k V_j^*)_{12}|^2 (x_{j,1} - x_{j,2})(x_{k,1} - x_{k,2}) \\ &= -\ell(\mathbf{x}_1) - \ell(\mathbf{x}_2) + \ell_S(\hat{X}, V), \end{aligned} \quad (5.15)$$

where  $\ell(\cdot)$  is defined in the first line of (5.6) and  $\ell_S(\hat{X}, V)$  is defined as

$$\ell_S(\hat{X}, V) = \frac{1}{2} \sum_{j,k} \mathfrak{s}_{jk} |(V_k V_j^*)_{12}|^2 (x_{j,1} - x_{j,2})(x_{k,1} - x_{k,2}). \quad (5.16)$$

Analogously to the notation  $L(D_{\pm}, I)$ , we will use  $K(D_{\pm}, I)$  to represent the value of  $K(\hat{X}, V)$  at  $(\hat{X}_j, V_j) = (D_{\pm}, I)$  for all  $j = 1, \dots, W$ . In addition,  $K(D_+, I)$  and  $K(D_-, I)$  are defined in the same manner. Observing that

$$\ell_S(D_{\pm}, I) = \ell_S(D_+, I) = \ell_S(D_-, I) = 0, \quad (5.17)$$

we have

$$K(D_{\pm}, I) = -\ell(a_+) - \ell(a_-), \quad K(D_+, I) = -2\ell(a_+), \quad K(D_-, I) = -2\ell(a_-). \quad (5.18)$$

Moreover, we employ the notation

$$\hat{K}(\hat{X}, V) = K(\hat{X}, V) - K(D_{\pm}, I) = -\ell_{++}(\mathbf{x}_1) - \ell_{+-}(\mathbf{x}_2) + \ell_S(\hat{X}, V). \quad (5.19)$$

We will need the following elementary observations that are easy to check from (5.18) and (5.6)

$$K(D_{\pm}, I) + L(D_{\pm}, I) = 0, \quad \operatorname{Re} K(D_+, I) = \operatorname{Re} K(D_-, I) = \operatorname{Re} K(D_{\pm}, I). \quad (5.20)$$

In addition, we introduce the  $W \times W$  matrix

$$S^v = (\mathfrak{s}_{jk}^v), \quad \mathfrak{s}_{jk}^v := \mathfrak{s}_{jk} |(V_k V_j^*)_{12}|^2, \quad (5.21)$$

and the  $2W \times 2W$  matrices

$$\mathbb{S} = S \oplus S, \quad \mathbb{S}^v := \mathbb{S} + \begin{pmatrix} -S^v & S^v \\ S^v & -S^v \end{pmatrix}, \quad (5.22)$$

where  $\mathbb{S}^v$  depends on the  $V$ -variables according to (5.21). Here we regard  $V$ -variables as fixed parameters. Due to the fact  $|(V_k V_j^*)_{12}| \in \mathbb{I}$ , it is easy to see that  $\mathbb{S}^v$  is a weighted Laplacian of a graph with  $2W$  vertices. In particular,  $\mathbb{S}^v \leq 0$ . By the definition (5.21), one can see that  $S_{ii}^v = 0$  for all  $i = 1, \dots, W$ . Consequently, we can obtain

$$\sum_{k \neq j} \mathbb{S}_{jk}^v = \sum_{k \neq j} \mathbb{S}_{kj}^v = -\mathbb{S}_{jj}^v = \begin{cases} -\mathfrak{s}_{jj}, & \text{if } j = 1, \dots, W \\ -\mathfrak{s}_{j-W, j-W}, & \text{if } j = W+1, \dots, 2W \end{cases}$$

Similarly to (1.4), we get

$$I + \mathbb{S}^v \geq c_0 I, \quad (5.23)$$

where  $c_0$  is the constant in Assumption 1.1 (ii). Moreover, it is not difficult to see from the definitions in (5.16), (5.21) and (5.22) that

$$\frac{1}{4} \sum_{j,k} \mathfrak{s}_{jk} \text{Tr}(\hat{X}_j - \hat{X}_k)^2 + \ell_S(\hat{X}, V) = -\frac{1}{2} \mathbf{x}' \mathbb{S}^v \mathbf{x}, \quad (5.24)$$

where we used the notation  $\mathbf{x} := (\mathbf{x}'_1, \mathbf{x}'_2)'$ .

Now let

$$\vartheta_j = \arg x_{j,1}, \quad \vartheta_{W+j} = \arg x_{j,2}, \quad \forall j = 1, \dots, W. \quad (5.25)$$

Then, recalling the parametrization of  $V_j$ 's in (3.24), we have the following lemma.

**Lemma 5.3.** *Assume that  $x_{j,1}, x_{j,2} \in \Sigma$  for all  $j = 1, \dots, W$ . We have*

$$\text{Re} \hat{K}(\hat{X}, V) \geq \frac{1}{4} \sum_{j,k=1}^{2W} (\mathbb{S}^v)_{jk} (\cos \vartheta_j - \cos \vartheta_k)^2 + c \sum_{j=1}^{2W} \left( \sin \vartheta_j - \frac{E}{2} \right)^2 \quad (5.26)$$

for some positive constant  $c$ . In addition,  $\text{Re} \hat{K}(\hat{X}, V)$  attains its minimum 0 at the following three types of saddle points

- Type I : For each  $j$ ,  $\hat{X}_j = D_{\pm}$  or  $D_{\mp}$ ,  
 $\theta_j \in \mathbb{L}$   $v_j = 0$  if  $\hat{X}_j = \hat{X}_1$ , and  $v_j = 1$  if  $\hat{X}_j \neq \hat{X}_1$ ,
- Type II : For each  $j$ ,  $\hat{X}_j = D_+$ ,  $V_j \in \mathring{U}(2)$ ,
- Type III : For each  $j$ ,  $\hat{X}_j = D_-$ ,  $V_j \in \mathring{U}(2)$ ,

which are the restrictions of three types of saddle points in Section 4.1, on  $\hat{X}$  and  $V$ -variables.

**Remark 5.4.** The Type I saddle points of  $(\hat{X}, V)$  are exactly those points satisfying

$$V_j^* \hat{X}_j V_j = D_{\pm}, \quad \forall j = 1, \dots, W, \quad \text{or} \quad V_j^* \hat{X}_j V_j = D_{\mp}, \quad \forall j = 1, \dots, W.$$

In Lemma 5.3, we wrote them in terms of  $\hat{X}_j$ ,  $v_j$  and  $\theta_j$  in order to evoke the parameterization in (3.22) and (3.24).

*Proof.* By (5.15), (5.24), the definitions of the functions  $\ell(\cdot)$  in (5.7) and  $\mathbb{k}(\cdot)$  in (5.4), we can write

$$\hat{K}(\hat{X}, V) = -\frac{1}{2} \mathbf{x}' \mathbb{S}^v \mathbf{x} - \sum_{j=1}^W \left( \mathbb{k}(x_{j,1}) + \mathbb{k}(x_{j,2}) - \mathbb{k}(a_+) - \mathbb{k}(a_-) \right).$$

By using (5.25) and the fact  $|x_{j,a}| = 1$  for all  $j = 1, \dots, W$  and  $a = 1, 2$ , we can obtain via elementary calculation

$$\begin{aligned} \operatorname{Re} \mathring{K}(\hat{X}, V) &= \frac{1}{4} \sum_{j,k=1}^{2W} (\mathbb{S}^v)_{jk} (\cos \vartheta_j - \cos \vartheta_k)^2 \\ &\quad - \frac{1}{4} \sum_{j,k=1}^{2W} (\mathbb{S}^v)_{jk} (\sin \vartheta_j - \sin \vartheta_k)^2 + \sum_{j=1}^{2W} \left( \sin \vartheta_j - \frac{E}{2} \right)^2. \end{aligned} \quad (5.27)$$

In light of the fact  $\mathbb{S}^v \leq 0$  and (5.23), we have

$$I + \frac{1}{2} \mathbb{S}^v \geq I + \mathbb{S}^v \geq c_0 I. \quad (5.28)$$

Applying (5.28) to the last two terms on the r.h.s. of (5.27) yields (5.26).

Now, we show that  $\operatorname{Re} \mathring{K}(\hat{X}, V)$  attains its minimum 0 at three types of points listed in Lemma 5.3. It is elementary to check that these points are minimum points along the contour by plugging them into the definition of  $\operatorname{Re} \mathring{K}(\hat{X}, V)$ . In the sequel, we show that they are the only solutions to the equation

$$\operatorname{Re} \mathring{K}(\hat{X}, V) = 0. \quad (5.29)$$

At first, by the second term on the r.h.s. of (5.26), we see that for any solution to (5.29),

$$\sin \vartheta_j = \frac{E}{2}, \quad \forall j = 1, \dots, 2W,$$

which implies that  $x_{j,a} = a_+$  or  $a_-$  for all  $j = 1, \dots, W$  and  $a = 1, 2$ , by recalling the definition (5.25) and the definitions of  $a_+$  and  $a_-$  in Section 1.4. Consequently, for each  $j$ ,  $\hat{X}_j$  can only be one of  $D_\pm$ ,  $D_\mp$ ,  $D_+$  and  $D_-$ .

Suppose that  $\hat{X}_1 = D_+$ , we claim that  $\hat{X}_j = D_+$  for all  $j$ . Otherwise, owing to the fact that the graph  $\mathcal{G}$  is connected, there exists  $\{i, j\} \in \mathcal{E}$  such that  $\mathfrak{s}_{ij} > 0$  and  $\hat{X}_i = D_+$  but  $\hat{X}_j = D_\pm$ ,  $D_\mp$  or  $D_-$ . Without loss of generality, we assume  $\hat{X}_j = D_\pm$ . In this case, we use the fact

$$\begin{aligned} \operatorname{Re} \mathring{K}(\hat{X}, V) &\geq \frac{1}{4} (\mathbb{S}^v)_{ij} (\cos \vartheta_i - \cos \vartheta_j)^2 + \frac{1}{4} (\mathbb{S}^v)_{i+W,j} (\cos \vartheta_{i+W} - \cos \vartheta_j)^2 \\ &\quad + \frac{1}{4} (\mathbb{S}^v)_{i,j+W} (\cos \vartheta_i - \cos \vartheta_{j+W})^2 + \frac{1}{4} (\mathbb{S}^v)_{i+W,j+W} (\cos \vartheta_{i+W} - \cos \vartheta_{j+W})^2, \end{aligned} \quad (5.30)$$

which follows from (5.26) directly. Now, by the assumption  $\hat{X}_i = D_+$  while  $\hat{X}_j = D_\pm$ , we have

$$\cos \vartheta_i = \cos \vartheta_{i+W} = \cos \vartheta_j = \operatorname{Re}(a_+) = \frac{\sqrt{4-E^2}}{2}, \quad \cos \vartheta_{j+W} = \operatorname{Re}(a_-) = -\frac{\sqrt{4-E^2}}{2},$$

which together with (5.30) implies that

$$\operatorname{Re} \mathring{K}(\hat{X}, V) \geq \frac{4-E^2}{4} \left( (\mathbb{S}^v)_{i,j+W} + (\mathbb{S}^v)_{i+W,j+W} \right) = \frac{4-E^2}{4} \mathfrak{s}_{ij} > 0, \quad (5.31)$$

contradicting to (5.29). Analogously, we can show that  $\hat{X}_j$  can not be  $D_\mp$  or  $D_-$ . Consequently, for a solution to (5.29), if  $\hat{X}_1 = D_+$ , we have shown that  $\hat{X}_j = D_+$  for all  $j$ . Similarly, we can show that if  $\hat{X}_1 = D_-$ , then  $\hat{X}_j = D_-$  for all  $j$ . These two kinds of solutions are collected as the Type II and Type III saddle points, respectively.

What remains is to show that if  $\hat{X}_1 = D_\pm$  or  $D_\mp$ , the solution to (5.29) must be one of the Type I saddle points. We only show the case of  $\hat{X}_1 = D_\pm$ . Assume that  $\{1, i\} \in \mathcal{E}$  in the graph  $\mathcal{G}$ , i.e.  $\mathfrak{s}_{1i} > 0$ . At first, similarly to the discussion from (5.30) to (5.31), we can show that  $\hat{X}_i$  can only be  $D_\pm$  or  $D_\mp$ . If  $\hat{X}_i = D_\pm$ , then by using (5.30) with  $j = 1$ , we have

$$\operatorname{Re} \mathring{K}(\hat{X}, V) \geq \frac{4-E^2}{4} \left( (\mathbb{S}^v)_{i,1+W} + (\mathbb{S}^v)_{i+W,1} \right) = \frac{4-E^2}{2} \mathfrak{s}_{1i}^v \geq 0,$$

and the equality holds if and only if  $V_i = I$ , according to the assumption  $V_1 = I$  and the definition in (5.21). The discussion on the case of  $\hat{X}_i = D_{\mp}$  is analogous. Consequently, we have

$$V_i^* \hat{X}_i V_i = V_1^* \hat{X}_1 V_1 = D_{\pm}. \quad (5.32)$$

Since the graph  $\mathcal{G}$  is connected, we can show that (5.32) holds for all  $i = 1, \dots, W$ . Analogously, if  $\hat{X}_1 = D_{\mp}$ , we can show that  $V_j^* \hat{X}_j V_j = D_{\mp}$  for all  $j = 1, \dots, W$ . Therefore, we completed the proof of Lemma 5.3.  $\square$

**5.3. Vicinities of the saddle points.** Having studied the saddle points of  $L(\hat{B}, T)$  and  $K(\hat{X}, V)$ , we then introduce some small vicinities of them. To this end, we introduce the quantity

$$\Theta \equiv \Theta(N, \varepsilon_0) := WN^{\varepsilon_0} \quad (5.33)$$

for small positive constant  $\varepsilon_0$  which will be chosen later. Let  $\mathbf{a} = (a_1, \dots, a_W) \in \mathbb{C}^W$  be any complex vector. In the sequel, we adopt the notation

$$\begin{aligned} \mathbf{a} + d &:= (a_1 + d, \dots, a_W + d), \quad d\mathbf{a} := (da_1, \dots, da_W), \quad \forall d \in \mathbb{C}, \\ \arg(\mathbf{a}) &:= (\arg(a_1), \dots, \arg(a_W)). \end{aligned} \quad (5.34)$$

Now, we define the following domains .

$$\begin{aligned} \Upsilon_+^b &\equiv \Upsilon_+^b(N, \varepsilon_0) := \left\{ \mathbf{a} \in \Gamma^W : \|\mathbf{a} - a_+\|_2^2 \leq \frac{\Theta}{M} \right\}, \\ \Upsilon_-^b &\equiv \Upsilon_-^b(N, \varepsilon_0) := \left\{ \mathbf{a} \in \bar{\Gamma}^W : \|\mathbf{a} + a_-\|_2^2 \leq \frac{\Theta}{M} \right\}, \\ \Upsilon_+^x &\equiv \Upsilon_+^x(N, \varepsilon_0) := \left\{ \mathbf{a} \in \Sigma^W : \|\arg(a_+^{-1}\mathbf{a})\|_2^2 \leq \frac{\Theta}{M} \right\}, \\ \Upsilon_-^x &\equiv \Upsilon_-^x(N, \varepsilon_0) := \left\{ \mathbf{a} \in \Sigma^W : \|\arg(a_-^{-1}\mathbf{a})\|_2^2 \leq \frac{\Theta}{M} \right\}, \\ \Upsilon_S &\equiv \Upsilon_S(N, \varepsilon_0) := \left\{ \mathbf{a} \in \mathbb{R}_+^{W-1} : -\mathbf{a}' S^{(1)} \mathbf{a} \leq \frac{\Theta}{M} \right\}, \end{aligned} \quad (5.35)$$

where the superscripts  $b$  and  $x$  indicate that these will be domains of the corresponding variables. In order to define the vicinities of the Type I saddle points properly, we introduce the permutation  $\epsilon_j$  of  $\{1, 2\}$ , for each triple  $(x_{j,1}, x_{j,2}, v_j)$ . Specifically, recalling the fact of  $u_j = \sqrt{1 - v_j^2}$  from (3.24), we define

$$v_{j, \epsilon_j} \equiv v_{j, \epsilon_j}(\epsilon_1) := v_j \mathbf{1}(\epsilon_j = \epsilon_1) + u_j \mathbf{1}(\epsilon_j \neq \epsilon_1).$$

Denoting by  $\epsilon = (\epsilon_1, \dots, \epsilon_W)$  and  $\epsilon(a) = (\epsilon_1(a), \dots, \epsilon_W(a))$  for  $a = 1, 2$ , we set

$$\mathbf{x}_{\epsilon(a)} = (x_{1, \epsilon_1(a)}, \dots, x_{W, \epsilon_W(a)}), \quad a = 1, 2, \quad \mathbf{v}_{\epsilon} = (v_{2, \epsilon_2}, \dots, v_{W, \epsilon_W}). \quad (5.36)$$

With this notation, we now define the Type I, II, and III vicinities, parameterized by  $(\mathbf{b}_1, \mathbf{b}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{t}, \mathbf{v})$  of the corresponding saddle point types. We also define the special case of the Type I vicinity, namely, Type I' vicinity, corresponding to the Type I' saddle point defined in Section 4.1.

**Definition 5.5.**

- *Type I vicinity* :  $(\mathbf{b}_1, \mathbf{b}_2, \mathbf{x}_{\epsilon(1)}, \mathbf{x}_{\epsilon(2)}, \mathbf{t}, \mathbf{v}_{\epsilon}) \in \Upsilon_+^b \times \Upsilon_-^b \times \Upsilon_+^x \times \Upsilon_-^x \times \Upsilon_S \times \Upsilon_S$  for some  $\epsilon$ .
- *Type I' vicinity* :  $(\mathbf{b}_1, \mathbf{b}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{t}, \mathbf{v}) \in \Upsilon_+^b \times \Upsilon_-^b \times \Upsilon_+^x \times \Upsilon_-^x \times \Upsilon_S \times \Upsilon_S$ .
- *Type II vicinity* :  $(\mathbf{b}_1, \mathbf{b}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{t}, \mathbf{v}) \in \Upsilon_+^b \times \Upsilon_-^b \times \Upsilon_+^x \times \Upsilon_+^x \times \Upsilon_S \times \mathbb{I}^{W-1}$ .
- *Type III vicinity* :  $(\mathbf{b}_1, \mathbf{b}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{t}, \mathbf{v}) \in \Upsilon_+^b \times \Upsilon_-^b \times \Upsilon_-^x \times \Upsilon_-^x \times \Upsilon_S \times \mathbb{I}^{W-1}$ .

In the following discussion, the parameter  $\varepsilon_0$  in  $\Theta$  is allowed to be different from line to line. However, given  $\varepsilon_1$  in (1.7), we shall always choose  $\varepsilon_2$  in (1.15) and  $\varepsilon_0$  in (5.33) according to the rule

$$C\varepsilon_2 \leq \varepsilon_0 \leq \frac{\varepsilon_1}{C} \quad (5.37)$$

for some sufficiently large  $C > 0$ . Consequently, by Assumption 1.13 we have

$$N(\log N)^{-10} \geq M = M^{1-4\varepsilon_0} M^{4\varepsilon_0} \geq W^{(4+2\gamma+\varepsilon_1)(1-4\varepsilon_0)} M^{4\varepsilon_0} \gg W^{(4+2\gamma+4\varepsilon_0)} M^{4\varepsilon_0} = W^{2\gamma} \Theta^4. \quad (5.38)$$

To prove Theorem 1.14, we split the task into three steps. The first step is to exclude the integral outside the vicinities. Specifically, we will show the following lemma.

**Lemma 5.6.** *Under Assumptions 1.1 and 1.13, we have,*

$$\begin{aligned} \mathcal{I}(\Gamma^W, \bar{\Gamma}^W, \Sigma^W, \Sigma^W, \mathbb{R}_+^{W-1}, \mathbb{I}^{W-1}) &= 2^W \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \Upsilon_S) \\ &\quad + \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_+^x, \Upsilon_S, \mathbb{I}^{W-1}) \\ &\quad + \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_-^x, \Upsilon_-^x, \Upsilon_S, \mathbb{I}^{W-1}) + O(e^{-\Theta}). \end{aligned} \quad (5.39)$$

**Remark 5.7.** *The first three terms on the r.h.s. of (5.39) correspond to the integrals over vicinities of the Type I, II, and III saddle points, respectively. Note that for the first term, we have used the argument in Section 4.1, namely, the total contribution of the integral over the Type I vicinity is  $2^W$  times that over the Type I' vicinity.*

The second step, is to estimate the integral over the Type I vicinity. We have the following lemma.

**Lemma 5.8.** *Under Assumptions 1.1 and 1.13, there exists some positive constant  $C_0$  uniform in  $n$  and some positive number  $N_0 = N_0(n)$  such that for all  $N \geq N_0$ ,*

$$2^W |\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \Upsilon_S)| \leq \frac{N^{C_0}}{(N\eta)^n}. \quad (5.40)$$

The last step is to show that the integral over the Type II and III vicinities are also negligible.

**Lemma 5.9.** *Under Assumptions 1.1 and 1.13, there exists some positive constant  $c$  such that,*

$$\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_+^x, \Upsilon_S, \mathbb{I}^{W-1}) = O(e^{-cW}), \quad \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_-^x, \Upsilon_-^x, \Upsilon_S, \mathbb{I}^{W-1}) = O(e^{-cW}).$$

Therefore, the remaining task is to prove Lemmas 5.1, 5.6, 5.8 and 5.9. For the convenience of the reader, we outline the organization of the subsequent part as follows.

At first, the proofs of Lemmas 5.1 and 5.6 require a discussion on the bound of the integrand, especially on the term  $A(\cdot)$ , which contains the integral over all the Grassmann variables. To this end, we perform a crude analysis for the function  $A(\cdot)$  in Section 6 in advance, with which we are able to prove Lemmas 5.1 and 5.6 in Section 7.

Then, we can restrict ourselves to the integral over the vicinities, i.e., prove Lemmas 5.8 and 5.9. It will be shown that in the vicinity  $(\mathbf{b}_1, \mathbf{b}_2, \mathbf{t}) \in \Upsilon_+^b \times \Upsilon_-^b \times \Upsilon_S$ , the factor  $\exp\{-M\hat{L}(\hat{B}, T)\}$  is approximately the product of a complex Gaussian measure of the  $\hat{B}$ -variables and a real Gaussian measure of the  $\mathbf{t}$ -variables. Here, by “complex Gaussian measure” we mean a function of the form  $\exp\{-\mathbf{u}'\mathbf{A}\mathbf{u}\}$ , where  $\mathbf{u}$  is a real vector, while  $\mathbf{A}$  is a complex matrix with positive-definite Hermitian part. In order to estimate the integral against this Gaussian measure (in an approximate sense), we shall get rid of the  $o(1)$  term in the integral of the form

$$\int d\mathbf{u} \exp\{-\mathbf{u}'\mathbf{A}\mathbf{u} + o(1)\} \mathbf{f}(\mathbf{u}), \quad (5.41)$$

for some function  $\mathbf{f}$ , which however cannot be done directly, owing to the fact that  $\mathbf{A}$  is complex. In our case, this problem can be solved by further deforming the contours of the  $\hat{B}$ -variables, following the



steepest descent paths exactly in the vicinity. By doing this, we can get a real Gaussian measure, thus the remainder terms can be easily controlled when integrate against this measure. The situation for  $\exp\{-M\hat{K}(\hat{X}, V)\}$  is a little bit more complicated due to different types of the saddle points. However, in the Type I vicinity, we can do the same thing. Hence, in Section 8, we will analyze the Gaussian measure (in an approximate sense)  $\exp\{-M(\hat{K}(\hat{X}, V) + \hat{L}(\hat{B}, T))\}$ , especially, we will further deform the contours of  $\hat{X}$  and  $\hat{B}$ -variables in the vicinities, whereby we can prove Lemmas 5.8 in Section 9. In the Type II and III vicinities, we will bound  $\exp\{-M\hat{K}(\hat{X}, V)\}$  by its absolute value directly. It turns out to be enough for our proof of Lemma 5.9, which is given in Section 10.

## 6. CRUDE BOUND ON $A(\hat{X}, \hat{B}, V, T)$

In this section, we provide a bound on the function  $A(\cdot)$  in terms of the  $\hat{B}, T$ -variables, which holds on all the domains under discussion in the sequel. Here, by *crude bound* we mean a bound of order  $\exp\{O(WN^{\varepsilon_2})\}$ , which will be specified in Lemma 6.1 below. By the definition in (3.31), we see that  $A(\cdot)$  is an integral of the product of  $\mathcal{Q}(\cdot)$ ,  $\mathcal{P}(\cdot)$  and  $\mathcal{F}(\cdot)$ . We will mainly treat  $\mathcal{Q}(\cdot)$  as a function of  $\omega^{[1]}, \xi^{[1]}$ -variables, treat  $\mathcal{P}(\cdot)$  as a function of  $\Omega, \Xi$ -variables, and treat  $\mathcal{F}(\cdot)$  as a function of  $X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}, P_1, Q_1$ -variables. However, in the function  $\mathcal{Q}(\cdot)$ , we actually have every argument mentioned above. Hence, we perform the integral over  $\omega^{[1]}$ -variables and  $\xi^{[1]}$ -variables for  $\mathcal{Q}(\cdot)$  at first. The resulting function  $Q(\cdot)$  turns out to be a polynomial of the remaining arguments. As mentioned in Section 4.2, a typical procedure we will adopt is to ignore  $Q(\cdot)$  at first, then estimate the integrals of  $\mathcal{P}(\cdot)$  and  $\mathcal{F}(\cdot)$ , which are denoted by  $P(\cdot)$  and  $F(\cdot)$ , respectively (see (4.4) and (4.5)), finally, we make necessary comment on how to modify the bounding scheme to take  $Q(\cdot)$  into account, whereby we can get the desired bound for  $A(\cdot)$ .

For the sake of simplicity, from now on, we will use the notation

$$\begin{aligned} \omega_{j,1} &:= \omega_{j,11}, & \omega_{j,2} &:= \omega_{j,12}, & \omega_{j,3} &:= \omega_{j,21}, & \omega_{j,4} &:= \omega_{j,22}, \\ \xi_{j,1} &:= \xi_{j,11}, & \xi_{j,2} &:= \xi_{j,21}, & \xi_{j,3} &:= \xi_{j,12}, & \xi_{j,4} &:= \xi_{j,22}. \end{aligned} \quad (6.1)$$

Moreover, we introduce the domains

$$\begin{aligned} \hat{\Sigma} &:= \left\{ r e^{i\vartheta} : |r - 1| \leq \frac{1}{10}, \vartheta \in \mathbb{L} \right\}, \\ \mathbb{K} \equiv \mathbb{K}(E) &:= \begin{cases} \left\{ \omega \in \mathbb{C} : 0 \leq \arg \omega \leq \frac{\arg a_+}{2} + \frac{\pi}{8} \right\}, & \text{if } E \geq 0, \\ \left\{ \omega \in \mathbb{C} : \frac{\arg a_+}{2} - \frac{\pi}{8} \leq \arg \omega \leq 0 \right\}, & \text{if } E < 0. \end{cases} \end{aligned} \quad (6.2)$$

By the assumption that  $|E| \leq \sqrt{2} - \kappa$  in (1.15), it is easy to see that  $|\arg \omega| \leq \pi/4 - c$  for all  $\omega \in \mathbb{K} \cup \bar{\mathbb{K}}$ , where  $c$  is some positive constant depending on  $\kappa$ . Our aim is to show the following lemma.

**Lemma 6.1.** *Suppose that  $\mathbf{b}_1, \mathbf{b}_2, \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{K}^W \times \bar{\mathbb{K}}^W \times \hat{\Sigma}^W \times \hat{\Sigma}^W$ . Under the assumption of Theorem 1.14, we have*

$$|A(\hat{X}, \hat{B}, V, T)| \leq e^{O(WN^{\varepsilon_2})} \prod_{j=1}^W \left( r_{j,1}^{-1} + r_{j,2}^{-1} + t_j + 1 \right)^C := e^{O(WN^{\varepsilon_2})} \mathbf{p}(\mathbf{r}^{-1}, \mathbf{t}).$$

**Remark 6.2.** *Obviously, using the terminology introduced at the end of Section 3.6, we have*

$$\mathbf{p}(\mathbf{r}^{-1}, \mathbf{t}) \in \mathfrak{Q}\left(\{r_{j,1}^{-1}, r_{j,2}^{-1}, t_j\}_{j=1}^W; \kappa_1, \kappa_2, \kappa_3\right), \quad \kappa_1 = e^{O(W)}, \quad \kappa_2, \kappa_3 = O(1). \quad (6.3)$$

6.1. **Integral of  $\mathcal{Q}$ .** In this section, we investigate the function

$$\mathcal{Q}(\Omega, \Xi, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) := \int d\omega^{[1]} d\xi^{[1]} \mathcal{Q}(\Omega, \Xi, \omega^{[1]}, \xi^{[1]}, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}). \quad (6.4)$$

Recall  $\mathfrak{Q}_{\deg}(\mathbf{a}; \kappa_1, \kappa_2, \kappa_3)$  defined at the end of Section 3.6, the parameterization in (3.15) and (3.24) and the notation introduced in (6.1). We shall show the following lemma.

**Lemma 6.3.** *If we regard  $\sigma$ ,  $v_p^{[1]}$ ,  $v_q^{[1]}$ ,  $P_1$  and  $(X^{[1]})^{-1}$ -entries as fixed parameters, we have*

$$\mathcal{Q}(\cdot) \in \mathfrak{Q}_{\deg}(\mathfrak{S}; \kappa_1, \kappa_2, \kappa_3), \quad \kappa_1 = W^{O(1)}, \quad \kappa_2, \kappa_3 = O(1), \quad (6.5)$$

where  $\mathfrak{S}$  is the set of variables defined by

$$\mathfrak{S} := \left\{ t, s, (y_k^{[1]})^{-1}, e^{i\sigma_k^{[1]}}, e^{-i\sigma_k^{[1]}}, \frac{\omega_{i,a}\xi_{j,b}}{M} \right\}_{\substack{i,j=1,\dots,W; \\ k=p,q; a,b=1,\dots,4}}.$$

*Proof.* Note that  $\mathcal{Q}(\cdot)$  can be regarded as a function of the Grassmann variables in  $\omega^{[1]}$  and  $\xi^{[1]}$ . Hence, by the definition in (3.29), it is a polynomial of these variables with bounded degree. To have a closer look on this polynomial, we start with

$$\exp \left\{ -\frac{1}{\sqrt{M}} \sum_j \tilde{s}_{jk} Tr P_1^* \Omega_j Q_1 \mathbf{w}_k^{[1]} \xi_k^{[1]} J \right\}, \quad k = p, q. \quad (6.6)$$

Observe that, in the exponent of (6.6), one  $\xi_k^{[1]}$ -variable must be accompanied by one  $\Omega$ -variable. In addition, we combine the factor  $1/\sqrt{M}$  with  $\Omega_j$ 's. Then, by Taylor expansion with respect to  $\xi_k^{[1]}$ -variables, it is easy to find

$$\exp \left\{ -\frac{1}{\sqrt{M}} \sum_j \tilde{s}_{jk} Tr P_1^* \Omega_j Q_1 \mathbf{w}_k^{[1]} \xi_k^{[1]} J \right\} \in \mathfrak{Q}_{\deg}(\mathfrak{S}_{1,k}; \kappa_1, \kappa_2, \kappa_3), \quad \kappa_1 = W^{O(1)}, \quad \kappa_2, \kappa_3 = O(1), \quad (6.7)$$

where

$$\mathfrak{S}_{1,k} := \left\{ t, s, e^{i\sigma_k^{[1]}}, \frac{\omega_{j,a}\xi_{k,b}^{[1]}}{\sqrt{M}} \right\}_{\substack{j=1,\dots,W; \\ a=1,\dots,4; b=1,2}}, \quad k = p, q.$$

Analogously, we can show that for  $k = p, q$

$$\exp \left\{ -\frac{1}{\sqrt{M}} \sum_j \tilde{s}_{kj} Tr \mathbf{w}_k^{[1]} (\mathbf{w}_k^{[1]})^* J Q_1^{-1} \Xi_j P_1 \right\} \in \mathfrak{Q}_{\deg}(\mathfrak{S}_{2,k}; \kappa_1, \kappa_2, \kappa_3), \quad \kappa_1 = W^{O(1)}, \quad \kappa_2, \kappa_3 = O(1), \quad (6.8)$$

where

$$\mathfrak{S}_{2,k} := \left\{ t, s, e^{-i\sigma_k^{[1]}}, \frac{\xi_{j,a}\omega_{k,b}^{[1]}}{\sqrt{M}} \right\}_{\substack{j=1,\dots,W; \\ a=1,\dots,4; b=1,2}}, \quad k = p, q.$$

Here we have used the fact that  $Q_1^{-1}$ -entries are the same as  $Q_1$ -entries, up to a sign. In a similar manner, one can show that for  $k, \ell = p, q$ ,

$$\exp \left\{ -\frac{1}{M} \tilde{s}_{k\ell} Tr \mathbf{w}_k^{[1]} (\mathbf{w}_k^{[1]})^* J (\mathbf{w}_\ell^{[1]}) \xi_\ell^{[1]} J \right\} \in \mathfrak{Q}_{\deg}(\mathfrak{S}_{3,k,\ell}; \kappa_1, \kappa_2, \kappa_3), \quad \kappa_1, \kappa_2, \kappa_3 = O(1), \quad (6.9)$$

and

$$\left( 1 - (y_k^{[1]})^{-1} \xi_k^{[1]} (X_k^{[1]})^{-1} \omega_k^{[1]} \right)^2 \in \mathfrak{Q}_{\deg}(\mathfrak{S}_{4,k}; \kappa_1, \kappa_2, \kappa_3), \quad \kappa_1, \kappa_2, \kappa_3 = O(1), \quad (6.10)$$

where

$$\mathfrak{S}_{3,k,\ell} := \left\{ e^{i\sigma_k^{[1]}}, e^{-i\sigma_\ell^{[1]}}, \omega_{k,a}\xi_{\ell,b}^{[1]} \right\}_{a,b=1,2}, \quad \mathfrak{S}_{4,k} := \left\{ (y_k^{[1]})^{-1}, \omega_{k,a}\xi_{\ell,b}^{[1]} \right\}_{a,b=1,2}.$$

Hence, by (6.7)-(6.10) and (3.29), we see that

$$\mathcal{Q}(\cdot) \in \mathfrak{Q}_{\text{deg}}(\mathfrak{S}_5; \kappa_1, \kappa_2, \kappa_3), \quad \kappa_1 = W^{O(1)}, \quad \kappa_2, \kappa_3 = O(1), \quad (6.11)$$

where

$$\mathfrak{S}_5 := \left\{ t, s, (y_k^{[1]})^{-1}, e^{i\sigma_k^{[1]}}, e^{-i\sigma_k^{[1]}}, \frac{\omega_{j,r}\xi_{k,b}^{[1]}}{\sqrt{M}}, \frac{\xi_{j,r}\omega_{k,b}^{[1]}}{\sqrt{M}}, \omega_{k,a}^{[1]}\xi_{\ell,b}^{[1]} \right\}_{\substack{j=1,\dots,W; k=p,q; \\ r=1,\dots,4; a,b=1,2}}.$$

By the definition in (6.4),  $\mathcal{Q}(\cdot)$  is the integral of  $\mathcal{Q}(\cdot)$  over the  $\omega^{[1]}$  and  $\xi^{[1]}$ -variables. Now, we regard all the other variables in  $\mathfrak{S}_5$ , except  $\omega^{[1]}$  and  $\xi^{[1]}$ -variables, as parameters. By the definition of Grassmann integral, we know that only the coefficient of the highest order term  $\prod_{k=p,q} \prod_{a=1,2} \omega_{k,a}^{[1]} \xi_{k,a}^{[1]}$  in  $\mathcal{Q}(\cdot)$  survives after integrating  $\omega^{[1]}$  and  $\xi^{[1]}$ -variables out. Then, it is easy to see (6.5) from (6.11), completing the proof.  $\square$

**6.2. Integral of  $\mathcal{P}$ .** In this subsection, we temporarily ignore the  $\Omega$  and  $\Xi$ -variables from  $\mathcal{Q}(\cdot)$ , and estimate  $\mathcal{P}(\hat{X}, \hat{B}, V, T)$  defined in (4.4). Recalling  $r_{j,1}$  and  $r_{j,2}$  defined in (5.12), we can formulate our estimate as follows.

**Lemma 6.4.** *Suppose that the assumptions in Lemma 6.1 hold. We have*

$$|\mathcal{P}(\hat{X}, \hat{B}, V, T)| \leq e^{O(W)} \prod_{j=1}^W (r_{j,1}^{-1} + r_{j,2}^{-1} + t_j + 1)^{O(1)}. \quad (6.12)$$

*Proof.* We start with one factor from  $\mathcal{P}(\cdot)$  (see (3.28)), namely

$$\begin{aligned} \varpi_j &:= \frac{1}{\det^M (1 + M^{-1} V_j^* \hat{X}_j^{-1} V_j \Omega_j T_j^{-1} \hat{B}_j^{-1} T_j \Xi_j)} \\ &= \exp \left\{ -M \log \det (1 + M^{-1} V_j^* \hat{X}_j^{-1} V_j \Omega_j T_j^{-1} \hat{B}_j^{-1} T_j \Xi_j) \right\} \\ &= 1 + \sum_{\ell=1}^4 \frac{1}{M^{\ell-1}} \mathfrak{p}_\ell(\hat{X}_j, \hat{B}_j, V_j, T_j, \Omega_j, \Xi_j). \end{aligned} \quad (6.13)$$

Here  $\mathfrak{p}_\ell(\cdot)$  is a polynomial in  $\hat{X}_j^{-1}$ ,  $\hat{B}_j^{-1}$ ,  $V_j$ ,  $T_j$ ,  $\Omega_j$  and  $\Xi_j$ -entries with bounded degree and bounded coefficients. Here we used the fact that  $V_j^*$  and  $T_j^{-1}$ -entries are the same as  $V_j$  and  $T_j$ -entries, respectively, up to a sign. Moreover, if we regard  $\mathfrak{p}_\ell(\cdot)$  as a polynomial of  $\Omega_j$  and  $\Xi_j$ -entries, it is homogeneous, with degree  $2\ell$ , and the total degree for  $\Omega_j$ -variables is  $\ell$ , thus that for  $\Xi_j$ -entries is also  $\ell$ . More specifically, we can write

$$\mathfrak{p}_\ell(\hat{X}_j, \hat{B}_j, V_j, T_j, \Omega_j, \Xi_j) = \sum_{\substack{\alpha_1, \dots, \alpha_\ell, \\ \beta_1, \dots, \beta_\ell=1}}^4 \mathfrak{p}_{\ell, \alpha, \beta}(\hat{X}_j, \hat{B}_j, V_j, T_j) \prod_{i=1}^{\ell} \omega_{j, \alpha_i} \xi_{j, \beta_i},$$

where we used the notation in (6.1) and denoted  $\alpha = (\alpha_1, \dots, \alpha_\ell)$  and  $\beta = (\beta_1, \dots, \beta_\ell)$ . It is easy to verify that  $\varpi_j$  is of the form (6.13) by taking Taylor expansion with respect to the Grassmann variables. The expansion in (6.13) terminates at  $\ell = 4$ , owing to the fact that there are totally 8 Grassmann variables from  $\Omega_j$  and  $\Xi_j$ . In addition, it is also easy to check that  $\mathfrak{p}_{\ell, \alpha, \beta}(\cdot)$  is a polynomial of  $\hat{X}_j^{-1}$ ,  $\hat{B}_j^{-1}$ ,  $V_j$ ,  $T_j$ -entries with bounded degree and bounded coefficients, which implies that there exist two positive constants  $C_1$  and  $C_2$ , such that

$$|\mathfrak{p}_{\ell, \alpha, \beta}(\cdot)| \leq C_1 (r_{j,1}^{-1} + r_{j,2}^{-1} + t_j + 1)^{C_2} \quad (6.14)$$

uniformly in  $\ell$ ,  $\alpha$  and  $\beta$ . Here we used the fact that  $\hat{X}_j^{-1}$  and  $V_j$ -entries are all bounded and  $T_j$ -entries are bounded by  $1 + t_j$ .

Now, we go back to the definition of  $\mathcal{P}(\cdot)$  in (3.28) and study the last factor. Similarly to the discussion above, it is easy to see that for  $k = p$  or  $q$ ,

$$\begin{aligned}\hat{\omega}_k &:= \frac{\det(V_k^* \hat{X}_k V_k + M^{-1} \Omega_k T_k^{-1} \hat{B}_k^{-1} T_k \Xi_k)}{\det \hat{B}_k} \\ &= \hat{\mathbf{p}}_0(\hat{X}_k, \hat{B}_k) + \sum_{\ell=1}^4 \sum_{\alpha, \beta} \hat{\mathbf{p}}_{\ell, \alpha, \beta}(\hat{X}_k, \hat{B}_k, V_k, T_k) \prod_{i=1}^{\ell} \omega_{k, \alpha_{k,i}} \xi_{k, \beta_{k,i}},\end{aligned}\quad (6.15)$$

where  $\hat{\mathbf{p}}_0(\cdot) = \det \hat{X}_k / \det \hat{B}_k$  and  $\hat{\mathbf{p}}_{\ell, \alpha, \beta}(\cdot)$ 's are some polynomials of  $\hat{X}_k$ ,  $\hat{B}_k^{-1}$ ,  $V_k$ ,  $T_k$ -entries with bounded degree and bounded coefficients. Similarly, we have

$$|\hat{\mathbf{p}}_0(\cdot)|, |\hat{\mathbf{p}}_{\ell, \alpha, \beta}(\cdot)| \leq C_1 (r_{k,1}^{-1} + r_{k,2}^{-1} + t_k + 1)^{C_2} \quad (6.16)$$

for some positive constants  $C_1$  and  $C_2$ .

According to the definitions in (6.13) and (6.15), we can rewrite (3.28) as

$$\mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T) = \exp \left\{ - \sum_{j,k} \tilde{\mathbf{s}}_{jk} Tr \Omega_j \Xi_k \right\} \prod_{j=1}^W \varpi_j \prod_{k=p,q} \hat{\omega}_k. \quad (6.17)$$

In light of the discussion above,  $\prod_{j=1}^W \varpi_j \prod_{k=p,q} \hat{\omega}_k$  is a polynomial of  $\hat{X}^{-1}$ ,  $\hat{B}^{-1}$ ,  $V$ ,  $T$ ,  $\Omega$  and  $\Xi$ -entries, in which each monomial is of the form

$$\mathbf{q}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}}(\hat{X}^{-1}, \hat{B}^{-1}, V, T) \prod_{j=1}^W \prod_{i=1}^{\ell_j} \omega_{j, \alpha_{j,i}} \xi_{j, \beta_{j,i}}, \quad (6.18)$$

where we used the notation

$$\begin{aligned}\vec{\ell} &= (\ell_1, \dots, \ell_W), \quad \vec{\alpha} = (\alpha_1, \dots, \alpha_W), \quad \vec{\beta} = (\beta_1, \dots, \beta_W), \\ \alpha_j &= (\alpha_{j,1}, \dots, \alpha_{j, \ell_j}), \quad \beta_j = (\beta_{j,1}, \dots, \beta_{j, \ell_j}),\end{aligned}$$

and  $\mathbf{q}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}}(\cdot)$  is a polynomial of  $\hat{X}$ ,  $\hat{X}^{-1}$ ,  $\hat{B}^{-1}$ ,  $V$  and  $T$ -entries. Moreover, all the entries of  $\vec{\ell}$ ,  $\vec{\alpha}$  and  $\vec{\beta}$  are bounded by 4. By (6.14) and (6.16), we have

$$|\mathbf{q}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}}(\hat{X}^{-1}, \hat{B}^{-1}, V, T)| \leq e^{O(W)} \prod_{j=1}^W (r_{j,1}^{-1} + r_{j,2}^{-1} + t_j + 1)^C. \quad (6.19)$$

In addition, it is easy to see that the number of the summands of the form (6.18) in  $\prod_{j=1}^W \varpi_j \prod_{k=p,q} \hat{\omega}_k$  is bounded by  $e^{O(W)}$ .

Define the vectors

$$\vec{\Omega} := (\omega_1, \omega_2, \omega_3, \omega_4), \quad \vec{\Xi} := (\xi_1, \xi_2, \xi_3, \xi_4), \quad (6.20)$$

where

$$\omega_\alpha = (\omega_{1,\alpha}, \dots, \omega_{W,\alpha}), \quad \xi_\alpha = (\xi_{1,\alpha}, \dots, \xi_{W,\alpha}), \quad \alpha = 1, 2, 3, 4.$$

Here we used the notation (6.1). In addition, we introduce the matrix

$$\tilde{\mathbb{H}} = \tilde{S} \oplus \tilde{S} \oplus \tilde{S} \oplus \tilde{S}.$$

It is easy to check

$$\sum_{j,k} \tilde{\mathbf{s}}_{jk} Tr \Omega_j \Xi_k = \vec{\Omega} \tilde{\mathbb{H}} \vec{\Xi}'.$$

By using the Gaussian integral formula for the Grassmann variables (3.2), we see that for each  $\vec{\ell}$ ,  $\vec{\alpha}$  and  $\vec{\beta}$ , we have

$$\left| \int d\Omega d\Xi \cdot \exp \left\{ - \sum_{j,k} \tilde{s}_{jk} Tr \Omega_j \Xi_k \right\} \cdot \prod_{j=1}^W \prod_{i=1}^{\ell_j} \omega_{j,\alpha_{j,i}} \xi_{j,\beta_{j,i}} \right| \leq |\det \tilde{\mathbb{H}}^{(\mathbf{l}|\mathbf{J})}|, \quad (6.21)$$

for some index sets  $\mathbf{l}$  and  $\mathbf{J}$  with  $|\mathbf{l}| = |\mathbf{J}|$ . By Assumption 1.1 (i) and (ii), we see that the 2-norm of each row of  $\tilde{S}$  is  $O(1)$ . Consequently, by using Hadamard's inequality, we have

$$|\det \tilde{\mathbb{H}}^{(\mathbf{l}|\mathbf{J})}| = e^{O(W)}. \quad (6.22)$$

Therefore, (6.17)-(6.22) and the bound  $e^{O(W)}$  for the total number of summands of the form (6.18) in  $\prod_{j=1}^W \varpi_j \prod_{k=p,q} \hat{\omega}_k$  imply that

$$|\mathbb{P}(\hat{X}, \hat{B}, V, T)| \leq e^{O(W)} \prod_{j=1}^W (r_{j,1}^{-1} + r_{j,2}^{-1} + t_j + 1)^{O(1)}.$$

Thus we completed the proof.  $\square$

**6.3. Integral of  $\mathcal{F}$ .** In this subsection, we also temporarily ignore the  $X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}, P_1, Q_1$ -variables from  $\mathcal{Q}(\cdot)$ , and estimate  $\mathbb{F}(\hat{X}, \hat{B}, V, T)$  defined in (4.5). We have the following lemma.

**Lemma 6.5.** *Suppose that the assumptions in Lemma 6.1 hold. We have*

$$|\mathbb{F}(\hat{X}, \hat{B}, V, T)| \leq e^{O(WN^{\varepsilon_2})} \prod_{k=p,q} (r_{k,1}^{-1} + r_{k,2}^{-1} + t_k + 1)^{O(1)}.$$

*Proof.* Recalling the decomposition of  $\mathcal{F}(\cdot)$  in (3.18) together with the parameterization in (3.27), we will study the integrals

$$\mathbb{G}(\hat{B}, T) := \int d\nu(Q_1) d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} g(Q_1, T, \hat{B}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}), \quad (6.23)$$

$$\mathbb{F}(\hat{X}, V) := \int d\mu(P_1) dX^{[1]} f(P_1, V, \hat{X}, X^{[1]}) \quad (6.24)$$

separately. Recalling the convention at the end of Section 3, we use  $f(\cdot)$  and  $g(\cdot)$  to represent the integrands above. One can refer to (3.19) and (3.20) for the definition.

From the assumption  $\eta \leq M^{-1}N^{\varepsilon_2}$ , it is easy to see

$$|\mathbb{F}(\hat{X}, V)| \leq e^{O(WN^{\varepsilon_2})}, \quad (6.25)$$

since  $P_1, V, \hat{X}, X^{[1]}$ -variables are all bounded and  $|\det X_p^{[1]}|, |\det X_q^{[1]}| \sim 1$  when  $\mathbf{x}_1, \mathbf{x}_2 \in \hat{\Sigma}$  defined in (6.2).

For  $\mathbb{G}(\hat{B}, T)$ , we use the facts

$$\begin{aligned} \operatorname{Re}(Tr B_j Y_k^{[1]} J) &\geq 0, \quad \operatorname{Re}(i Tr Y_k^{[1]} J Z) = -\eta Tr Y_k^{[1]} \leq 0, \quad Tr Y_k^{[1]} J Y_\ell^{[1]} J \geq 0, \quad k, \ell = p, q, \\ |(\mathbf{w}_q^{[1]} (\mathbf{w}_q^{[1]})^*)_{12}| &\leq 1, \quad |(\mathbf{w}_p^{[1]} (\mathbf{w}_p^{[1]})^*)_{21}| \leq 1, \end{aligned} \quad (6.26)$$

to estimate trivially several terms, whereby we can get the bound

$$|g(\cdot)| \leq \exp \left\{ -M\eta \sum_{j=1}^W Tr \operatorname{Re}(B_j) J \right\} \prod_{k=p,q} (y_k^{[1]})^{n+3} \exp \left\{ -\tilde{s}_{kk} Tr \operatorname{Re}(B_k) Y_k^{[1]} J \right\}. \quad (6.27)$$

Here  $\operatorname{Re}(B_j) = Q_1^{-1} T_j^{-1} \operatorname{Re}(\hat{B}_j) T_j Q_1$ . Hence, integrating  $y_p^{[1]}$  and  $y_q^{[1]}$  out yields

$$\int_{\mathbb{R}_+^2} dy_p^{[1]} dy_q^{[1]} |g(Q_1, T, \hat{B}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]})| \leq C \frac{\exp \left\{ -M\eta \sum_{j=1}^W Tr \operatorname{Re}(B_j) J \right\}}{\prod_{k=p,q} \left( (\mathbf{w}_k^{[1]})^* J \operatorname{Re}(B_k) \mathbf{w}_k^{[1]} \right)^{C_1}}, \quad (6.28)$$

for some positive constants  $C$  and  $C_1$  depending on  $n$ , where we used the elementary facts that  $\tilde{s}_{kk} \geq c$  for some positive constant  $c$  and

$$\text{Tr}(\text{Re}(B_j)Y_k^{[1]}J = y_k^{[1]}(\mathbf{w}_k^{[1]})^*J\text{Re}(B_j)\mathbf{w}_k^{[1]}, \quad k = p, q, \quad j = 1, \dots, W. \quad (6.29)$$

Now, note that

$$(\mathbf{w}_k^{[1]})^*J\text{Re}B_j\mathbf{w}_k^{[1]} \geq \lambda_1(J\text{Re}B_j), \quad k = p, q, \quad j = 1, \dots, W. \quad (6.30)$$

In addition, it is also easy to see  $\lambda_1(T_j) = s_j - t_j$  and  $\lambda_1(Q_1) = s - t$ , according to the definitions in (3.24). Now, by the fact  $JA^{-1} = AJ$  for any  $A \in \dot{U}(1, 1)$ , we have

$$J\text{Re}B_j = Q_1T_j\text{diag}(\text{Re}b_{j,1}, \text{Re}b_{j,2})T_jQ_1. \quad (6.31)$$

Consequently, we can get

$$\lambda_1(J\text{Re}B_j) \geq (s_j - t_j)^2(s - t)^2 \min\{\text{Re}b_{j,1}, \text{Re}b_{j,2}\} = \frac{\min\{\text{Re}b_{j,1}, \text{Re}b_{j,2}\}}{(s_j + t_j)^2(s + t)^2}, \quad (6.32)$$

by recalling the facts  $s^2 - t^2 = 1$  and  $s_j^2 - t_j^2 = 1$ . Therefore, combining (6.28), (6.30) and (6.32), we have

$$\int_{\mathbb{R}_+^2} dy_p^{[1]} dy_q^{[1]} |g(\cdot)| \leq C(s + t)^{4C_1} \exp \left\{ -M\eta \sum_{j=1}^W \text{Tr}(\text{Re}B_j)J \right\} \prod_{k=p,q} \frac{(s_k + t_k)^{2C_1}}{(\min\{\text{Re}b_{k,1}, \text{Re}b_{k,2}\})^{C_1}}. \quad (6.33)$$

Now, what remains is to estimate the exponential function in (6.33). By elementary calculation from (6.31) we obtain

$$\text{Tr}(\text{Re}B_j)J \geq (\text{Re}b_{j,1} + \text{Re}b_{j,2})((s_j^2 + t_j^2)(s^2 + t^2) - 4sts_jt_j).$$

Observe that

$$(s_j^2 + t_j^2)(s^2 + t^2) - 4sts_jt_j = \frac{(s_j^2 + t_j^2)^2(s^2 + t^2)^2 - 16(sts_jt_j)^2}{(s_j^2 + t_j^2)(s^2 + t^2) + 4sts_jt_j} \geq \frac{4t^4 + 4t^2 + 4t_j^4 + 4t_j^2 + 1}{2(1 + 2t_j^2)(1 + 2t^2)} \geq \frac{1 + 2t^2}{2(1 + 2t_j^2)}.$$

It implies that

$$\begin{aligned} \exp \left\{ -M\eta \sum_{j=1}^W \text{Tr}(\text{Re}B_j)J \right\} &\leq \exp \left\{ -2M\eta \sum_{j=1}^W \frac{\text{Re}b_{j,1} + \text{Re}b_{j,2}}{1 + 2t_j^2} (1 + 2t^2) \right\} \\ &\leq \exp \left\{ -cM\eta \sum_{j=1}^W \frac{r_{j,1} + r_{j,2}}{1 + 2t_j^2} (1 + 2t^2) \right\}, \end{aligned} \quad (6.34)$$

for some positive constant  $c$ , where in the last step we used the fact

$$\text{Re}b_{j,\alpha} \geq cr_{j,\alpha}, \quad \forall j = 1, \dots, W, \alpha = 1, 2, \quad (6.35)$$

for some positive constant  $c$ , in light of the assumption  $|E| \leq \sqrt{2} - \kappa$  and the definition of  $\mathbb{K}$  in (6.2). Plugging (6.34) into (6.33), estimating  $(s + t)^2 \leq 2(1 + 2t^2)$ , and integrating  $t$  out, we can crudely bound

$$\int_{\mathbb{R}_+^2} dy_p^{[1]} dy_q^{[1]} \int_{\mathbb{R}^+} 2t dt \cdot |g(\cdot)| \leq C \left( \frac{1}{M\eta} \right)^{C_2} \left( \sum_{j=1}^W \frac{r_{j,1} + r_{j,2}}{1 + 2t_j^2} \right)^{-C_2} \prod_{k=p,q} \frac{(1 + 2t_k^2)^{C_1}}{((\min\{\text{Re}b_{k,1}, \text{Re}b_{k,2}\}))^{C_1}}. \quad (6.36)$$

Now, we use the trivial bounds

$$\left( \sum_{j=1}^W \frac{r_{j,1} + r_{j,2}}{1 + 2t_j^2} \right)^{-C_2} \leq \left( \frac{1 + 2t_p^2}{r_{p,1} + r_{p,2}} \right)^{C_2} \leq \left( (1 + 2t_p^2)(r_{p,1}^{-1} + r_{p,2}^{-1}) \right)^{C_2}, \quad (6.37)$$

and

$$\frac{1 + 2t_k^2}{\min\{\text{Re}b_{k,1}, \text{Re}b_{k,2}\}} \leq C(1 + 2t_k^2)(r_{k,1}^{-1} + r_{k,2}^{-1}). \quad (6.38)$$

Inserting (6.37) and (6.38) into (6.36) and integrating out the remaining variables yields

$$|\mathbb{G}(\hat{B}, T)| \leq C \left( \frac{1}{M\eta} \right)^{C_1} \prod_{k=p,q} \left( r_{k,1}^{-1} + r_{k,2}^{-1} + t_k + 1 \right)^{C_3}. \quad (6.39)$$

Combining (6.25) and (6.39) we can get the bound

$$|\mathbb{F}(\hat{X}, \hat{B}, V, T)| \leq e^{O(WN^{\varepsilon_2})} \prod_{k=p,q} \left( r_{k,1}^{-1} + r_{k,2}^{-1} + t_k + 1 \right)^{O(1)}.$$

Hence, we completed the proof of Lemma 6.5.  $\square$

**6.4. Summing up: Proof of Lemma 6.1.** In the discussions in Sections 6.2 and 6.3, we ignored the irrelevant factor  $\mathbb{Q}(\cdot)$ . However, it is easy to modify the discussion slightly to take this factor into account, whereby we can prove Lemma 6.1.

*Proof of Lemma 6.1.* At first, by the definition in (6.4), we can rewrite (3.31) as

$$\mathbb{A}(\cdot) = \int dX^{[1]} d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\Omega d\Xi d\mu(P_1) d\nu(Q_1) \mathcal{P}(\cdot) \mathbb{Q}(\cdot) \mathcal{F}(\cdot).$$

Now, by the conclusion  $\kappa_1 = W^{O(1)}$  in Lemma 6.3, it suffices to consider one term in  $\mathbb{Q}(\cdot)$ , which is a monomial of the form

$$\mathbf{p}(t, s, (y_p^{[1]})^{-1}, (y_q^{[1]})^{-1}) \mathbf{q}(\Omega, \Xi),$$

regarding  $\sigma, v_p^{[1]}, v_q^{[1]}, P_1$ -variables,  $X^{[1]}$ -variables and  $\mathbf{w}^{[1]}$ -variables as bounded parameters. Here  $\mathbf{p}(\cdot)$  is a monomial of  $t, s, (y_p^{[1]})^{-1}, (y_q^{[1]})^{-1}$  and  $\mathbf{q}(\cdot)$  is a monomial of  $\Omega, \Xi$ -variables, both with bounded coefficients and bounded degrees, according to the fact  $\kappa_2, \kappa_3 = O(1)$  in Lemma 6.3. Now we define

$$\begin{aligned} \mathbb{P}_q(\hat{X}, \hat{B}, V, T) &:= \int d\Omega d\Xi \mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T) \cdot \mathbf{q}(\Omega, \Xi), \\ \mathbb{F}_p(\hat{X}, \hat{B}, V, T) &:= \int dX^{[1]} d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\mu(P_1) d\nu(Q_1) \\ &\quad \times \mathcal{F}(\hat{X}, \hat{B}, V, T, P_1, Q_1, X^{[1]}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \cdot \mathbf{p}(t, s, (y_p^{[1]})^{-1}, (y_q^{[1]})^{-1}). \end{aligned} \quad (6.40)$$

By repeating the discussions in Sections 6.2 and 6.3 with slight modification, we can easily see that

$$\begin{aligned} |\mathbb{P}_q(\hat{X}, \hat{B}, V, T)| &\leq e^{O(W)} \prod_{j=1}^W (r_{j,1}^{-1} + r_{j,2}^{-1} + t_j + 1)^{O(1)}, \\ |\mathbb{F}_p(\hat{X}, \hat{B}, V, T)| &\leq e^{O(WN^{\varepsilon_2})} \prod_{k=p,q} \left( r_{k,1}^{-1} + r_{k,2}^{-1} + t_k + 1 \right)^{O(1)} \end{aligned}$$

hold as well. Therefore, we completed the proof of Lemma 6.1.  $\square$

## 7. PROOFS OF LEMMAS 5.1 AND 5.6

In this section, with the aid of Lemma 6.1, we prove Lemmas 5.1 and 5.6. According to Lemmas 5.2 and 5.3, one can see that away from the saddles,  $\text{Re}\hat{L}(\hat{B}, T)$  and  $\text{Re}\hat{K}(\hat{X}, V)$  increase quadratically in  $\hat{B}$ -variables and  $\hat{X}$ -variables, respectively. Hence, it is easy to control the integral (5.2) over these variables outside the vicinities. However, like  $\hat{B}$ -variables, the domain of  $\mathbf{t}$ -variables is also not compact. This forces us to analyze the exponential function

$$\mathbb{M}(\mathbf{t}) := \exp \left\{ -M \text{Re}(\ell_S(\hat{B}, T)) \right\} \quad (7.1)$$

carefully for any fixed  $\hat{B}$ -variables.

Recall the definition of the sector  $\mathbb{K}$  in (6.2). For  $\mathbf{b}_1 \in \mathbb{K}^W$  and  $\mathbf{b}_2 \in \bar{\mathbb{K}}^W$ , we have

$$\min_{j,k} \operatorname{Re}(b_{j,1} + b_{j,2})(b_{k,1} + b_{k,2}) \geq c \min_{j,k} \sum_{a,b=1,2} r_{j,a} r_{k,b} \geq c \min_{j,a} r_{j,a}^2 := \mathfrak{A}(\hat{B}), \quad (7.2)$$

for some positive constant  $c$  depending on  $\kappa$  from (1.15). From now on, we regard  $\mathbb{M}(\mathbf{t})$  as a measure of the  $\mathbf{t}$ -variables and study it in the following two regions separately:

$$(i) : \mathbf{t} \in \mathbb{I}^{W-1}, \quad (ii) : \mathbf{t} \in \mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}.$$

Roughly speaking, when  $\mathbf{t} \in \mathbb{I}^{W-1}$ , we will see that  $\mathbb{M}(\mathbf{t})$  can be bounded pointwisely by a Gaussian measure. More specifically, we have the following lemma.

**Lemma 7.1.** *With the notation above, we have*

$$\mathbb{M}(\mathbf{t}) \leq \exp \left\{ -\frac{M}{12} \mathfrak{A}(\hat{B}) \sum_{j,k} \mathfrak{s}_{jk} (t_k - t_j)^2 \right\}, \quad \forall \mathbf{t} \in \mathbb{I}^{W-1}.$$

However, the behavior of  $\mathbb{M}(\mathbf{t})$  for  $\mathbf{t} \in \mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}$  is much more sophisticated. We will not try to provide a pointwise control of  $\mathbb{M}(\mathbf{t})$  in this region. Instead, we will bound the integral of  $\mathbf{q}(\mathbf{t})$  against  $\mathbb{M}(\mathbf{t})$  over this region, for any given monomial  $\mathbf{q}(\cdot)$  of interest. More specifically, recalling the definition of  $\Theta$  in (5.33) and the spanning tree  $\mathcal{G}_0 = (\mathcal{V}, \mathcal{E}_0)$  in Assumption 1.1, and additionally setting

$$\mathfrak{L} := \frac{M}{4} \mathfrak{A}(\hat{B}) \min_{i,j \in \mathcal{E}_0} \mathfrak{s}_{ij}, \quad (7.3)$$

we have the following lemma.

**Lemma 7.2.** *Let  $\mathbf{q}(\mathbf{t}) = \prod_{j=2}^W t_j^{n_j}$  be a monomial of  $\mathbf{t}$ -variables, with powers  $n_j = O(1)$  for all  $j = 2, \dots, W$ . We have*

$$\int_{\mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}} \prod_{j=2}^W dt_j \mathbb{M}(\mathbf{t}) \mathbf{q}(\mathbf{t}) \leq \left(1 + \mathfrak{L}^{-\frac{1}{2}}\right)^{O(W^2)} \exp \left\{ -\Theta^2 \mathfrak{A}(\hat{B}) + O(W^2 \log N) \right\} \quad (7.4)$$

**Remark 7.3.** *Roughly speaking, by Lemma 7.2 we see that the integral of  $\mathbf{q}(\mathbf{t})$ -variables against the measure  $\mathbb{M}(\mathbf{t})$  over the region  $\mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}$  is exponentially small, owing to the fact  $\Theta^2 \gg W^2 \log N$ .*

We will postpone the proofs of Lemmas 7.1 and 7.2 to the end of this section. In the sequel, at first, we prove Lemmas 5.1 and 5.6 with the aid of Lemmas 6.1, 7.1 and 7.2. Before commencing the formal proofs, we mention two basic facts which are formulated as the following lemma.

**Lemma 7.4.** *Under Assumption 1.1, we have the following two facts.*

- For the smallest eigenvalue of  $S^{(1)}$ , there exists some positive constant  $c$  such that

$$\lambda_1(-S^{(1)}) \geq \frac{c}{W^2}. \quad (7.5)$$

- Let  $\boldsymbol{\varrho} = (\varrho_2, \dots, \varrho_W)'$  be a real vector and  $\varrho_1 = 0$ . If there is at least one  $\alpha \in \{2, \dots, W\}$  such that  $\varrho_\alpha \geq \Theta/\sqrt{M}$ , then we have

$$\sum_{j,k} \mathfrak{s}_{jk} (\varrho_j - \varrho_k)^2 \geq \frac{\Theta}{M}. \quad (7.6)$$

*Proof.* Let  $\boldsymbol{\varrho} = (\varrho_2, \dots, \varrho_W)'$  be a real vector and  $\varrho_1 = 0$ . Now, we assume  $|\varrho_\alpha| = \max_{\beta=2, \dots, W} |\varrho_\beta|$ . Then

$$\frac{-\boldsymbol{\rho}' S^{(1)} \boldsymbol{\rho}}{\|\boldsymbol{\rho}\|_2^2} = \frac{\frac{1}{2} \sum_{j,k} \mathfrak{s}_{jk} (\rho_j - \rho_k)^2}{\sum_j \rho_j^2} \geq \frac{c(\rho_\alpha - \rho_1)^2}{W^2 \rho_\alpha^2} = \frac{c}{W^2},$$



where the second step follows from Assumption 1.1 (iv) and Cauchy-Schwarz inequality. Analogously, we have

$$\sum_{j,k} \mathfrak{s}_{jk} (\varrho_j - \varrho_k)^2 \geq \frac{c}{W} \varrho_\alpha^2 \geq \frac{\Theta}{M}$$

according to the definition of  $\Theta$  in (5.33). Hence, we completed the proof.  $\square$

Recalling the notation defined in (5.2) and the facts  $|x_{j,a}| = 1$  and  $|b_{j,a}| = r_{j,a}$  for all  $j = 1, \dots, W$  and  $a = 1, 2$ , for any sequence of domains, we have

$$\begin{aligned} & |\mathcal{I}(\mathbf{I}_1^b, \mathbf{I}_2^b, \mathbf{I}_1^x, \mathbf{I}_2^x, \mathbf{I}^t, \mathbf{I}^v)| \\ & \leq e^{O(W \log N)} \int_{\mathbf{I}_1^b} \prod_{j=1}^W db_{j,1} \int_{\mathbf{I}_2^b} \prod_{j=1}^W db_{j,2} \int_{\mathbf{I}_1^x} \prod_{j=1}^W dx_{j,1} \int_{\mathbf{I}_2^x} \prod_{j=1}^W dx_{j,2} \int_{\mathbf{I}^t} \prod_{j=2}^W 2t_j dt_j \int_{\mathbf{I}^v} \prod_{j=2}^W 2v_j dv_j \\ & \quad \times \exp \left\{ -M(\operatorname{Re} K(\hat{X}, V) + \operatorname{Re} L(\hat{B}, T)) \right\} \cdot |A(\hat{X}, \hat{B}, V, T)| \cdot \prod_{j=1}^W (r_{j,1} + r_{j,2})^2. \end{aligned} \quad (7.7)$$

In addition, according to Lemma 6.1, we have

$$|A(\hat{X}, \hat{B}, V, T)| \cdot \prod_{j=2}^W 2t_j \cdot \prod_{j=2}^W 2v_j \cdot \prod_{j=1}^W (r_{j,1} + r_{j,2})^2 \leq e^{O(WN^{\varepsilon_2})} \tilde{\mathbf{p}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}), \quad (7.8)$$

for some polynomial  $\tilde{\mathbf{p}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t})$  with positive coefficients, and

$$\tilde{\mathbf{p}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}) \in \mathfrak{Q} \left( \{r_{j,1}, r_{j,2}, r_{j,1}^{-1}, r_{j,2}^{-1}, t_j\}_{j=1}^W; \kappa_1, \kappa_2, \kappa_3 \right), \quad \kappa_1 = e^{O(W)}, \quad \kappa_2, \kappa_3 = O(1). \quad (7.9)$$

**7.1. Proof of Lemma 5.1.** At first, since throughout the whole proof, the domains of  $\mathbf{x}_1$ ,  $\mathbf{x}_2$ , and  $\mathbf{v}$ -variables, namely,  $\Sigma^W$ ,  $\Sigma^W$  and  $\mathbb{I}^{W-1}$ , will not be involved, we just use  $*$ 's to represent them, in order to simplify the notation.

Now, we introduce the following contours with the parameter  $\mathfrak{D} \in \mathbb{R}_+$ ,

$$\Gamma_{\mathfrak{D}} := \{ra_+ | r \in [0, \mathfrak{D}]\} \subset \Gamma, \quad \mathbb{R}_{\mathfrak{D}} = [0, (\operatorname{Re} a_+) \mathfrak{D}] \subset \mathbb{R}_+, \quad \mathcal{L}_{\mathfrak{D}} := \{(\operatorname{Re} a_+) \mathfrak{D} + i(\operatorname{Im} a_+)r | r \in [0, \mathfrak{D}]\}.$$

In addition, we recall the sector  $\mathbb{K}$  defined in (6.2). Then, trivially, we have

$$\mathbb{R}_+, \Gamma, \mathcal{L}_{\mathfrak{D}} \subset \mathbb{K}, \quad \mathbb{R}_+, \bar{\Gamma}, \bar{\mathcal{L}}_{\mathfrak{D}} \in \bar{\mathbb{K}}, \quad \forall \mathfrak{D} \in \mathbb{R}_+.$$

We claim that the integrand in (5.2) is an analytic function of the  $\hat{B}$ -variables. To see this, we can go back to the integral representation (3.17) and the definitions of  $L(B)$  and  $\mathcal{P}(\Omega, \Xi, X, B)$  in (3.18). Note that since  $\exp\{M \log \det B_j\} = (\det B_j)^M$ , actually the logarithmic terms in  $L(B)$  do not produce any singularity in the integrand in (3.17). In addition, according to the fact that the  $\chi^\ell = 0$  for any Grassmann variable  $\chi$  and  $\ell \geq 2$ , the factors  $\det^M(1 + M^{-1} X_j^{-1} \Omega_j B_j^{-1} \Xi_j)$  is actually a polynomial of  $\Omega_j$ ,  $\Xi_j$ ,  $X_j^{-1}$  and  $B_j^{-1}$ -entries with degree 16. The other factors containing  $\hat{B}$ -variables  $\mathcal{P}(\cdot)$  can be checked analogously. Hence, it is easy to see that  $\exp\{-ML(\cdot)\} \mathcal{P}(\cdot)$  is analytic in  $\hat{B}$ -variables. Consequently, we have

$$\mathcal{I} \left( (\Gamma_{\mathfrak{D}} \cup \mathcal{L}_{\mathfrak{D}})^W, (\bar{\Gamma}_{\mathfrak{D}} \cup \bar{\mathcal{L}}_{\mathfrak{D}})^W, *, *, \mathbb{R}_+^{W-1}, * \right) = \mathcal{I} \left( (\mathbb{R}_{\mathfrak{D}})^W, (\mathbb{R}_{\mathfrak{D}})^W, *, *, \mathbb{R}_+^{W-1}, * \right).$$

Hence, to prove Lemma 5.1, it suffices to prove the following lemma.

**Lemma 7.5.** *Suppose that  $|E| \leq \sqrt{2} - \kappa$ . As  $\mathfrak{D} \rightarrow \infty$ , the following convergence hold,*

$$\begin{aligned} (i) : \quad & \mathcal{I} \left( (\Gamma_{\mathfrak{D}} \cup \mathcal{L}_{\mathfrak{D}})^W, (\bar{\Gamma}_{\mathfrak{D}} \cup \bar{\mathcal{L}}_{\mathfrak{D}})^W, *, *, \mathbb{R}_+^{W-1}, * \right) - \mathcal{I} \left( (\Gamma_{\mathfrak{D}})^W, (\bar{\Gamma}_{\mathfrak{D}})^W, *, *, \mathbb{R}_+^{W-1}, * \right) \rightarrow 0, \\ (ii) : \quad & \mathcal{I} \left( (\Gamma)^W, (\bar{\Gamma})^W, *, *, \mathbb{R}_+^{W-1}, * \right) - \mathcal{I} \left( (\Gamma_{\mathfrak{D}})^W, (\bar{\Gamma}_{\mathfrak{D}})^W, *, *, \mathbb{R}_+^{W-1}, * \right) \rightarrow 0, \end{aligned}$$

$$(iii): \quad \mathcal{I}(\mathbb{R}_+^W, \mathbb{R}_+^W, *, *, \mathbb{R}_+^{W-1}, *) - \mathcal{I}(\mathbb{R}_{\mathfrak{D}}^W, \mathbb{R}_{\mathfrak{D}}^W, *, *, \mathbb{R}_+^{W-1}, *) \rightarrow 0.$$

*Proof.* For simplicity, we use the notation

$$\begin{aligned} \mathbf{I}_{\mathfrak{D}}^{b,1} &:= (\Gamma_{\mathfrak{D}} \cup \mathcal{L}_{\mathfrak{D}})^W \times (\bar{\Gamma}_{\mathfrak{D}} \cup \bar{\mathcal{L}}_{\mathfrak{D}})^W \setminus (\Gamma_{\mathfrak{D}})^W \times (\bar{\Gamma}_{\mathfrak{D}})^W, \\ \mathbf{I}_{\mathfrak{D}}^{b,2} &:= \Gamma^W \times \bar{\Gamma}^W \setminus (\Gamma_{\mathfrak{D}})^W \times (\bar{\Gamma}_{\mathfrak{D}})^W, \\ \mathbf{I}_{\mathfrak{D}}^{b,3} &:= \mathbb{R}_+^W \times \mathbb{R}_+^W \setminus \mathbb{R}_{\mathfrak{D}}^W \times \mathbb{R}_{\mathfrak{D}}^W. \end{aligned}$$

Now, recall the definition of the function  $\ell(\mathbf{a})$  in (5.6) and the representation of  $L(\hat{B}, T)$  in (5.5). Hence, in light of the definition of  $\mathbb{M}(\mathbf{t})$  in (7.1), we have

$$\exp \{ -M \operatorname{Re} L(\hat{B}, T) \} = \exp \{ -M (\operatorname{Re} \ell(\mathbf{b}_1) + \operatorname{Re} \ell(-\mathbf{b}_2)) \} \mathbb{M}(\mathbf{t}). \quad (7.10)$$

By the assumption  $|E| \leq \sqrt{2} - \kappa$ , we see that  $\operatorname{Re} b_{j,a} b_{k,a} > 0$  for all  $b_{j,a}, b_{k,a} \in \mathbb{K} \cup \bar{\mathbb{K}}$ . Consequently, when  $b_{j,1} \in \mathbb{K}$  and  $b_{j,2} \in \bar{\mathbb{K}}$  for all  $j = 1, \dots, W$ , we have for some positive constant  $c$  dependent of  $\kappa$  in (1.15),

$$\begin{aligned} \operatorname{Re} \ell(\mathbf{b}_1) + \operatorname{Re} \ell(-\mathbf{b}_2) &\geq \sum_{a=1,2} \sum_j \left( \frac{1}{2} (1 + \mathfrak{s}_{jj}) \operatorname{Re} b_{j,a}^2 + (-1)^{a+1} E \operatorname{Im} b_{j,a} - \log r_{j,a} \right) \\ &\geq c \sum_{a=1,2} \sum_j r_{j,a}^2 - \sum_{a=1,2} \sum_j \log r_{j,a}, \end{aligned} \quad (7.11)$$

where we used Assumption 1.1 (ii) and the fact that  $(-1)^{a+1} E \operatorname{Im} b_{j,a} \geq 0$ .

Now, when  $(\mathbf{b}_1, \mathbf{b}_2) \in \mathbf{I}_{\mathfrak{D}}^{b,i}$  for  $i = 1, 2, 3$ , we have  $\sum_{a=1,2} \sum_j r_{j,a}^2 \geq c \mathfrak{D}^2$  for some positive constant  $c$ , which implies the trivial fact

$$\sum_{a=1,2} \sum_j r_{j,a}^2 \geq \frac{1}{2} \sum_{a=1,2} \sum_j r_{j,a}^2 + \frac{c}{2} \mathfrak{D}^2. \quad (7.12)$$

Consequently, we can get from (7.10), (7.11) and (7.12) that for some positive constant  $c$ ,

$$\exp \{ -M \operatorname{Re} L(\hat{B}, T) \} \leq e^{-cM\mathfrak{D}^2} \prod_{a=1,2} \prod_{j=1}^W e^{-cMr_{j,a}^2} r_{j,a}^M \cdot \mathbb{M}(\mathbf{t})$$

holds in  $\mathbf{I}_{\mathfrak{D}}^{b,i}$  for  $i = 1, 2, 3$ . In addition, by the boundedness of  $V$  and  $\hat{X}$ -variables, we can get the trivial bound  $MK(\hat{X}, V) = O(N)$ . Hence, from (7.7) and (7.8) we see that the quantities in Lemma 7.5 (i), (ii) and (iii) can be bounded by the following integral with  $i = 1, 2, 3$ , respectively,

$$\begin{aligned} &e^{O(W \log N)} \int_{\mathbf{I}_{\mathfrak{D}}^{b,i}} \prod_{j=1}^W db_{j,1} db_{j,2} \int_{\Sigma^W} \prod_{j=1}^W dx_{j,1} \int_{\Sigma^W} \prod_{j=1}^W dx_{j,2} \int_{\mathbb{R}_+^{W-1}} \prod_{j=2}^W 2t_j dt_j \int_{\mathbb{I}^{W-1}} \prod_{j=2}^W 2v_j dv_j \\ &\times \exp \left\{ -M (\operatorname{Re} K(\hat{X}, V) + \operatorname{Re} L(\hat{B}, T)) \right\} \cdot |\mathbf{A}(\hat{X}, \hat{B}, V, T)| \cdot \prod_{j=1}^W (r_{j,1} + r_{j,2})^2 \\ &\leq e^{O(N)} e^{-cM\mathfrak{D}^2} \int_{\mathbf{I}_{\mathfrak{D}}^{b,i}} \prod_{j=1}^W db_{j,1} db_{j,2} \int_{\mathbb{R}_+^{W-1}} \prod_{j=2}^W dt_j \prod_{a=1,2} \prod_{j=1}^W e^{-cMr_{j,a}^2} r_{j,a}^M \cdot \mathbb{M}(\mathbf{t}) \cdot \tilde{\mathbf{p}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}). \end{aligned} \quad (7.13)$$

According to the facts  $\kappa_1 = e^{O(W)}$  and  $\kappa_2 = O(1)$  in (7.9), it suffices to consider one monomial in  $\tilde{\mathbf{p}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t})$  with bounded coefficient. That means, it suffices to estimate the integral

$$\int_{\mathbf{I}_{\mathfrak{D}}^{b,i}} \prod_{j=1}^W db_{j,1} db_{j,2} \int_{\mathbb{R}_+^{W-1}} dt_j \prod_{a=1,2} \prod_{j=1}^W e^{-cMr_{j,a}^2} r_{j,a}^M \cdot \mathbb{M}(\mathbf{t}) \cdot \tilde{\mathbf{q}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}), \quad i = 1, 2, 3, \quad (7.14)$$

for some monomial

$$\tilde{\mathbf{q}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}) = \prod_{a=1,2} \prod_{j=1}^W r_{j,a}^{\ell_{j,a}} \prod_{j=2}^W t_j^{n_j}, \quad |\ell_{j,a}|, n_j = O(1), \quad \text{for } j = 1, \dots, W, \quad a = 1, 2, \quad (7.15)$$

where the bound on  $\ell_j$ 's and  $n_j$ 's follows from the fact that  $\kappa_3 = O(1)$  in (7.9).

Bounding  $t_j$ 's by 1 trivially in the region  $\mathbf{t} \in \mathbb{I}^{W-1}$  and using Lemma 7.2 in the region  $\mathbf{t} \in \mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}$ , we can get for  $i = 1, 2, 3$ ,

$$(7.14) \leq e^{O(W^2 \log N)} \int_{\mathbf{I}_{\mathfrak{D}}^{b,i}} \prod_{j=1}^W db_{j,1} db_{j,2} \cdot \prod_{a=1,2} \prod_{j=1}^W e^{-cM r_{j,a}^2} r_{j,a}^{M+\ell_{j,a}} (1 + \mathfrak{L}^{-\frac{1}{2}})^{O(W^2)}.$$

By the definition of  $\mathfrak{A}(\hat{B})$  in (7.2) and the assumption  $M \gg W^4$  we see that

$$\prod_{a=1,2} \prod_{j=1}^W r_{j,a}^{M+\ell_{j,a}} (1 + \mathfrak{L}^{-\frac{1}{2}})^{O(W^2)} \leq \prod_{a=1,2} \prod_{j=1}^W r_{j,a}^{M(1+o(1))}.$$

Consequently, by using elementary Gaussian integral, we can get the trivial bound

$$(7.14) \leq e^{O(N \log N)}, \quad i = 1, 2, 3,$$

and then we have for  $i = 1, 2, 3$ ,

$$(7.13) \leq e^{-cM\mathfrak{D}^2 + O(N \log N)} \rightarrow 0, \quad \text{as } \mathfrak{D} \rightarrow \infty.$$

Thus we completed the proof.  $\square$

**7.2. Proof of Lemma 5.6.** Plugging the first identity of (5.20) and (7.8) into (7.7), we can write

$$\begin{aligned} |\mathcal{I}(\mathbf{I}_1^b, \mathbf{I}_2^b, \mathbf{I}_1^x, \mathbf{I}_2^x, \mathbf{I}^t, \mathbf{I}^v)| &\leq e^{O(WN^{\varepsilon_2})} \int_{\mathbf{I}_1^b} \prod_{j=1}^W db_{j,1} \int_{\mathbf{I}_2^b} \prod_{j=1}^W db_{j,2} \int_{\mathbf{I}_1^x} \prod_{j=1}^W dx_{j,1} \int_{\mathbf{I}_2^x} \prod_{j=1}^W dx_{j,2} \\ &\quad \times \int_{\mathbf{I}^t} \prod_{j=2}^W dt_j \int_{\mathbf{I}^v} \prod_{j=2}^W dv_j \exp \left\{ -M \left( \text{Re} \dot{K}(\hat{X}, V) + \text{Re} \dot{L}(\hat{B}, T) \right) \right\} \tilde{\mathbf{p}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}), \end{aligned} \quad (7.16)$$

where  $\tilde{\mathbf{p}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t})$  is specified in (7.9).

Lemma 5.6 immediately follows from the following two lemmas.

**Lemma 7.6.** *Under Assumptions 1.1 and 1.13, we have*

$$\mathcal{I}(\Gamma^W, \bar{\Gamma}^W, \Sigma^W, \Sigma^W, \mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}, \mathbb{I}^{W-1}) \leq e^{-\Theta^2}. \quad (7.17)$$

**Lemma 7.7.** *Under Assumptions 1.1 and 1.13, we have*

$$\begin{aligned} \mathcal{I}(\Gamma^W, \bar{\Gamma}^W, \Sigma^W, \Sigma^W, \mathbb{I}^{W-1}, \mathbb{I}^{W-1}) &= 2^W \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \Upsilon_S) \\ &\quad + \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \mathbb{I}^{W-1}) \\ &\quad + \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \mathbb{I}^{W-1}) + O(e^{-\Theta}). \end{aligned} \quad (7.18)$$

In the sequel, we prove Lemmas 7.6 and 7.7.

*Proof of Lemma 7.6.* Recall (7.16) with the choice of the integration domains

$$(\mathbf{I}_1^b, \mathbf{I}_2^b, \mathbf{I}_1^x, \mathbf{I}_2^x, \mathbf{I}^t, \mathbf{I}^v) = (\Gamma^W, \bar{\Gamma}^W, \Sigma^W, \Sigma^W, \mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}, \mathbb{I}^{W-1}).$$

To simplify the integral on the r.h.s. of (7.16), we use the fact  $\text{Re} \dot{K}(\hat{X}, V) \geq 0$  implied by (5.26), together with the facts that the  $\mathbf{x}$  and  $\mathbf{v}$ -variables are bounded by 1. Consequently, we can eliminate the integral over  $\mathbf{x}$  and  $\mathbf{v}$ -variables from the integral on the r.h.s. of (7.16). Moreover, according to (7.9), it suffices to prove

$$\int_{\Gamma^W} \prod_{j=1}^W db_{j,1} \int_{\bar{\Gamma}^W} \prod_{j=1}^W db_{j,2} \int_{\mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}} \prod_{j=2}^W dt_j \exp \left\{ -M \text{Re} \dot{L}(\hat{B}, T) \right\} \cdot \tilde{\mathbf{q}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}) \leq e^{-\Theta^2} \quad (7.19)$$

instead, where  $\tilde{\mathbf{q}}(\cdot)$  is the monomial defined in (7.15).

Now, by the first inequality of (5.13), we have

$$\exp \left\{ -M \operatorname{Re} \hat{L}(\hat{B}, T) \right\} \leq \exp \left\{ -M \sum_{a=1,2} \sum_{j=1}^W \left( c(r_{j,a} - 1)^2 + (r_{j,a} - \log r_{j,a} - 1) \right) \right\} \cdot \mathbb{M}(\mathbf{t}). \quad (7.20)$$

At first, we integrate  $\mathbf{t}$ -variables out by using Lemma 7.2, namely,

$$\begin{aligned} & \int_{\mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}} \prod_{j=2}^W dt_j \mathbb{M}(\mathbf{t}) \cdot \tilde{\mathbf{q}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}) \\ & \leq \prod_{a=1,2} \prod_{j=1}^W r_{j,a}^{\ell_{j,a}} \cdot \left( 1 + \mathfrak{L}^{-\frac{1}{2}} \right)^{O(W^2)} \exp \left\{ -\Theta^2 \mathfrak{A}(\hat{B}) + O(W^2 \log N) \right\}. \end{aligned} \quad (7.21)$$

Recall the definitions of  $\mathfrak{L}$  and  $\mathfrak{A}(\hat{B})$  in (7.3) and (7.2), and the fact  $W^2 = o(M^{1/2})$ . The proof of the following fact is an elementary exercise

$$\sup_{r \in \mathbb{R}_+} \exp \left\{ -M(r - \log r - 1) \right\} r^\ell = O(1), \quad \text{if } \ell = O(M^{\frac{1}{2}}). \quad (7.22)$$

Hence, we get the bound

$$\begin{aligned} & \prod_{a=1,2} \prod_{j=1}^W \exp \left\{ -M(r_{j,a} - \log r_{j,a} - 1) \right\} \cdot r_{j,a}^{\ell_{j,a}} \cdot \left( 1 + \mathfrak{L}^{-\frac{1}{2}} \right)^{O(W^2)} \\ & \leq \max_{b=1,2} \max_{k=1,\dots,W} \prod_{a=1,2} \prod_{j=1}^W \exp \left\{ -M(r_{j,a} - \log r_{j,a} - 1) \right\} \cdot r_{j,a}^{\ell_{j,a}} \left( 1 + r_{k,b}^{-1} \right)^{O(W^2)} = e^{O(W^2)}. \end{aligned} \quad (7.23)$$

Consequently, (7.20)-(7.23) imply that

$$\begin{aligned} & \int_{\mathbb{R}_+^{W-1} \setminus \mathbb{I}^{W-1}} \prod_{j=2}^W dt_j \exp \left\{ -M \operatorname{Re} \hat{L}(\hat{B}, T) \right\} \cdot \tilde{\mathbf{q}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}) \\ & \leq \exp \left\{ -\Theta^2 \mathfrak{A}(\hat{B}) + O(W^2 \log N) \right\} \cdot \prod_{a=1,2} \prod_{j=1}^W \exp \left\{ -cM(r_{j,a} - 1)^2 \right\} \\ & \leq e^{-\Theta^2} \cdot \prod_{a=1,2} \prod_{j=1}^W \exp \left\{ -\frac{c}{2} M(r_{j,a} - 1)^2 \right\}, \end{aligned} \quad (7.24)$$

for some positive constant  $c$ , where in the last step we use the obvious fact

$$\Theta^2 \mathfrak{A}(\hat{B}) + \frac{c}{2} \sum_{a=1,2} \sum_{j=1}^W M(r_{j,a} - 1)^2 \geq \Theta^2 \gg W^2 \log N$$

by (5.38) and the definition of  $\mathfrak{A}(\hat{B})$  in (7.2). Plugging the bound (7.24) into the l.h.s of (7.19) and taking the integral over  $\hat{B}$ -variables we can see that (7.19) holds, which further implies (7.17). Therefore, we completed the proof of Lemma 7.6.  $\square$

To prove Lemma 7.7, we split the exponential function into two parts. We use one part to control the integral, and the other will be estimated by its magnitude. More specifically, we shall prove the following two lemmas.

**Lemma 7.8.** *Under Assumptions 1.1 and 1.13, we have*

$$\begin{aligned} & \int_{\Gamma^W} \prod_{j=1}^W db_{j,1} \int_{\Gamma^W} \prod_{j=1}^W db_{j,2} \int_{\Sigma^W} \prod_{j=1}^W dx_{j,1} \int_{\Sigma^W} \prod_{j=1}^W dx_{j,2} \int_{\mathbb{I}^{W-1}} \prod_{j=2}^W dt_j \int_{\mathbb{I}^{W-1}} \prod_{j=2}^W dv_j \\ & \times \exp \left\{ -\frac{1}{2} M(\operatorname{Re} \hat{K}(\hat{X}, V) + \operatorname{Re} \hat{L}(\hat{B}, T)) \right\} \cdot \tilde{\mathbf{p}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}) \leq e^{O(W)}. \end{aligned} \quad (7.25)$$

**Lemma 7.9.** *If  $(\mathbf{b}_1, \mathbf{b}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{t}, \mathbf{v}) \in \Gamma^W \times \bar{\Gamma}^W \times \Sigma^W \times \Sigma^W \times \mathbb{I}^{W-1} \times \mathbb{I}^{W-1}$ , but not in any of the Types I, II, III vicinities in Definition 5.5, we have*

$$\exp \left\{ -\frac{1}{2} M(\text{Re}\hat{K}(\hat{X}, V) + \text{Re}\hat{L}(\hat{B}, T)) \right\} \leq e^{-\Theta}. \quad (7.26)$$

With Lemmas 7.8 and 7.9, we can prove Lemma 7.7.

*Proof of Lemma 7.7.* For the sake of simplicity, in this proof, we temporarily use  $\mathcal{I}_{\text{full}}$  to represent the l.h.s. of (7.18), i.e. the integral over the full domain, and use  $\mathcal{I}_I$ ,  $\mathcal{I}_{II}$  and  $\mathcal{I}_{III}$  to represent the first three terms on the r.h.s. of (7.18). Now, combining (7.16), (7.26) and (7.25), we see that,

$$|\mathcal{I}_{\text{full}} - \mathcal{I}_I - \mathcal{I}_{II} - \mathcal{I}_{III}| \leq e^{O(WN^{\varepsilon_2})} \cdot e^{-\Theta} \cdot e^{O(W)} \leq e^{-\Theta},$$

in light of the definition of  $\Theta$  in (5.33) and the assumption (5.37). Hence, we completed the proof of Lemma 7.7.  $\square$

*Proof of Lemma 7.8.* At first, again, the polynomial  $\tilde{\mathbf{p}}(\cdot)$  in the integrand can be replaced by the monomial  $\tilde{\mathbf{q}}(\cdot)$  defined in (7.15) in the discussion below, owing to the fact that  $\kappa_1 = \exp\{O(W)\}$  in (7.9). Then, the proof is similar to that of Lemma 7.6, but much simpler, since  $\mathbf{t}$ -variables are bounded by 1 now. Consequently, we can eliminate  $\hat{X}$ ,  $\mathbf{t}$  and  $\mathbf{v}$ -variables from the integral directly and use the trivial bounds

$$\tilde{\mathbf{q}}(\mathbf{r}, \mathbf{r}^{-1}, \mathbf{t}) \leq \prod_{a=1,2} \prod_{j=1}^W r_{j,a}^{\ell_{j,a}}, \quad \text{Re}\hat{L}(\hat{B}, T) \geq c \sum_{a=1,2} \sum_{j=1}^W (r_{j,a} - 1)^2, \quad (7.27)$$

where the latter is from (5.13). Hence, it suffices to show

$$\int_{\Gamma^W} \prod_{j=1}^W db_{j,1} \int_{\bar{\Gamma}^W} \prod_{j=1}^W db_{j,2} \prod_{a=1,2} \prod_{j=1}^W \exp \{ -cM(r_{j,a} - 1)^2 \} r_{j,a}^{\ell_{j,a}} \leq e^{O(W)}. \quad (7.28)$$

Note that (7.28) follows from elementary Gaussian integral immediately. Therefore, we completed the proof of Lemma 7.8.  $\square$

*Proof of Lemma 7.9.* At first, according to (5.13) and (5.26), we see both  $M\text{Re}\hat{L}(\hat{B}, T)$  and  $M\text{Re}\hat{K}(\hat{X}, V)$  are nonnegative on the full domain. Hence, it suffices to show one of them is larger than  $\Theta$  outside the Type I, II, III vicinities.

Note that for each type of vicinity, we have

$$(\mathbf{b}_1, \mathbf{b}_2, \mathbf{t}) \in \Upsilon_+^b \times \Upsilon_-^b \times \Upsilon_S. \quad (7.29)$$

Now, if (7.29) is violated, we have  $(\mathbf{b}_1, \mathbf{b}_2) \in \Gamma^W \times \bar{\Gamma}^W \setminus \Upsilon_+^b \times \Upsilon_-^b$  or  $\mathbf{t} \in \mathbb{I}^{W-1} \setminus \Upsilon_S$ . If the former holds, by using (5.13) and the definition of  $\Upsilon_+^b$  and  $\Upsilon_-^b$  in (5.35), we have

$$M\text{Re}\hat{L}(\hat{B}, T) \geq cM \sum_{a=1,2} \sum_{j=1}^W (r_{j,a} - 1)^2 = cM \|\mathbf{b}_1 - a_+\|_2^2 + cM \|\mathbf{b}_2 + a_-\|_2^2 \geq \Theta,$$

which shows (7.26) if  $(\mathbf{b}_1, \mathbf{b}_2) \in \Gamma^W \times \bar{\Gamma}^W \setminus \Upsilon_+^b \times \Upsilon_-^b$ .

Hence, it suffices to consider the case  $(\mathbf{b}_1, \mathbf{b}_2) \in \Upsilon_+^b \times \Upsilon_-^b$ ,  $\mathbf{t} \in \mathbb{I}^{W-1} \setminus \Upsilon_S$ . Using Lemma 7.1, we can see that

$$\mathbb{M}(\mathbf{t}) \leq \exp \left\{ -\frac{M}{6} \mathfrak{A}(\hat{B})(-\mathbf{t}S^{(1)}\mathbf{t}) \right\} \leq \exp \left\{ -\mathfrak{A}(\hat{B})\Theta \right\} \leq e^{-\Theta}, \quad (7.30)$$

where in the second step we used the definition of  $\Upsilon_S$  in (5.35) and in the last step we used the fact  $\mathfrak{A}(\hat{B}) \geq c$  if  $(\mathbf{b}_1, \mathbf{b}_2) \in \Upsilon_+^b \times \Upsilon_-^b$ . Then (7.20) and (7.30) also imply (7.26). Now, we turn to show

$M\text{Re}\hat{K}(\hat{X}, V) \geq \Theta$  outside the vicinities. Recalling the definition of  $\vartheta_j$ 's in (5.25), we split the discussion into two cases

$$(i) : \left( \sin \vartheta_j - \frac{E}{2} \right)^2 \leq \frac{\Theta}{M}, \quad \forall j = 1, \dots, 2W, \quad (ii) : \left( \sin \vartheta_j - \frac{E}{2} \right)^2 > \frac{\Theta}{M}, \quad \text{for some } j \in \{1, \dots, 2W\}.$$

Using (5.26), we can get  $M\text{Re}\hat{K}(\hat{X}, V) \geq \Theta$  in case (ii) immediately, so we can assume case (i) below. Then (i) implies that

$$|\arg(a_+^{-1}x_{j,a})|^2 \wedge |\arg(a_-^{-1}x_{j,a})|^2 \leq \frac{\Theta}{M}, \quad \forall j = 1, \dots, W; a = 1, 2. \quad (7.31)$$

Now, we claim that it suffices to focus on the following three subcases of (7.31),

(i') There exists a sequence of permutations of  $\{1, 2\}$ , namely,  $\epsilon = (\epsilon_1, \dots, \epsilon_W)$ , such that

$$|\arg(a_+^{-1}x_{j,\epsilon_j(1)})|^2 \leq \frac{\Theta}{M}, \quad |\arg(a_-^{-1}x_{j,\epsilon_j(2)})|^2 \leq \frac{\Theta}{M}, \quad \forall j = 1, \dots, W.$$

(ii') There exists

$$|\arg(a_+^{-1}x_{j,a})|^2 \leq \frac{\Theta}{M}, \quad \forall j = 1, \dots, W; a = 1, 2.$$

(iii') There exists

$$|\arg(a_-^{-1}x_{j,a})|^2 \leq \frac{\Theta}{M}, \quad \forall j = 1, \dots, W; a = 1, 2.$$

To see this, note that for those  $\hat{X}$ -variables which satisfy (7.31) but do not belong to any of the case (i'), (ii') or (iii') listed above, there must be a pair  $\{i, j\} \in \mathcal{E}$  such that

$$|\arg(a_+^{-1}x_{i,1})|^2, |\arg(a_+^{-1}x_{i,2})|^2 \leq \frac{\Theta}{M}, \quad |\arg(a_-^{-1}x_{j,1})|^2, |\arg(a_-^{-1}x_{j,2})|^2 \leq \frac{\Theta}{M}, \quad (7.32)$$

or there exists a permutation  $\epsilon_i$  such that

$$|\arg(a_+^{-1}x_{i,\epsilon_i(1)})|^2, |\arg(a_-^{-1}x_{i,\epsilon_i(2)})|^2 \leq \frac{\Theta}{M}, \quad |\arg(a_+^{-1}x_{j,1})|^2, |\arg(a_+^{-1}x_{j,2})|^2 \leq \frac{\Theta}{M}, \quad (7.33)$$

or

$$|\arg(a_+^{-1}x_{i,\epsilon_i(1)})|^2, |\arg(a_-^{-1}x_{i,\epsilon_i(2)})|^2 \leq \frac{\Theta}{M}, \quad |\arg(a_-^{-1}x_{j,1})|^2, |\arg(a_-^{-1}x_{j,2})|^2 \leq \frac{\Theta}{M}. \quad (7.34)$$

For each of (7.32), (7.33) and (7.34), we can perform a discussion similar to (5.30)-(5.31), to show that  $M\text{Re}\hat{K}(\hat{X}, V) \geq cM \gg \Theta$ .

Hence, it suffices to focus on cases (i'), (ii') and (iii') in the sequel. Now, we denote the domains of  $\hat{X}$ -variables satisfying (i'), (ii') and (iii') by  $\Upsilon_I^x$ ,  $\Upsilon_{II}^x$  and  $\Upsilon_{III}^x$ , respectively. In addition, in the remaining part of this proof, we temporarily use the term Type A vicinity to represent its restriction on  $\hat{X}$ -variables, for  $A = I, II, III$ . Obviously, we have

$$\text{Type A vicinity} \subset \Upsilon_A^x, \quad A = I, II, III.$$

If  $(\mathbf{x}_1, \mathbf{x}_2) \in \Upsilon_I^x$  but outside the Type I vicinity,

$$\| \arg(a_+^{-1}\mathbf{x}_{\epsilon(1)}) \|_2^2 + \| \arg(a_-^{-1}\mathbf{x}_{\epsilon(2)}) \|_2^2 \geq \frac{\Theta}{M}, \quad \forall \epsilon,$$

which easily implies that

$$M\text{Re}\hat{K}(\hat{X}, V) \geq cM \sum_{j=1}^{2W} \left( \sin \vartheta_j - \frac{E}{2} \right)^2 \geq \Theta. \quad (7.35)$$

Now, we turn to the case that there exists one sequence of permutations  $\epsilon$  such that  $(\mathbf{x}_{\epsilon(1)}, \mathbf{x}_{\epsilon(2)}) \in \Upsilon_+^x \times \Upsilon_-^x$ , but  $\mathbf{v}_\epsilon \notin \Upsilon_S$ . In this case, we just go back to the first line of (5.15) and do the transform

$$V_j \rightarrow \mathfrak{I}V_j, \quad \hat{X}_j \rightarrow \mathfrak{I}\hat{X}_j\mathfrak{I}$$

for those  $j$  with  $\epsilon(j) \neq \epsilon(1)$ , where  $\mathfrak{J}$  is defined in (4.1). Then, it suffices to consider

$$(\mathbf{x}_1, \mathbf{x}_2) \in \Upsilon_+^x \times \Upsilon_-^x, \quad \text{but} \quad \mathbf{v} \notin \Upsilon_S, \quad \text{or} \quad (\mathbf{x}_1, \mathbf{x}_2) \in \Upsilon_-^x \times \Upsilon_+^x, \quad \text{but} \quad \mathbf{v} \notin \Upsilon_S.$$

In either case, we can show that  $M\text{Re}\hat{K}(\hat{X}, V) \geq \Theta$ , analogously to case of  $(\mathbf{b}_1, \mathbf{b}_2) \in \Upsilon_+^x \times \Upsilon_-^x$  but  $\mathbf{t} \notin \Upsilon_S$ , in (7.30).

Now, what remains is to show that for those  $(\mathbf{x}_1, \mathbf{x}_2) \in \Upsilon_A^x$  but outside the Type A vicinity ( $A=II, III$ ), we have  $M\text{Re}\hat{K}(\hat{X}, V) \geq \Theta$ . We only discuss the case  $A = II$ , the other is analogous. Note that outside the Type II vicinity of  $\hat{X}$  variables we have

$$\|\arg(a_+^{-1}\mathbf{x}_1)\|_2^2 + \|\arg(a_+^{-1}\mathbf{x}_2)\|_2^2 \geq \frac{\Theta}{M}. \quad (7.36)$$

Observe that now we are already in  $\Upsilon_{II}^x$ , which means that all  $x_{j,a}$ 's are close to  $a_+$  and far away from  $a_-$ . That means, we have  $\sin(\arg(x_{ja})) - E/2 \sim \arg(a_+^{-1}x_{ja})$ . Consequently, (7.36) also implies (7.35). Therefore, we completed the proof of Lemma 7.9.  $\square$

**7.3. Proof of Lemma 7.1.** Using the definition of  $\ell_S(\hat{B}, T)$  in (5.6) and  $\mathfrak{A}(\hat{B})$  in (7.2) and the fact  $|(T_k T_j^{-1})_{12}| = |s_j t_k e^{i\sigma_k} - s_k t_j e^{i\sigma_j}|$ , we have

$$\text{Re}\ell_S(\hat{B}, T) \geq \frac{1}{2} \mathfrak{A}(\hat{B}) \sum_{j,k} s_{jk} |s_j t_k e^{i\sigma_k} - s_k t_j e^{i\sigma_j}|^2. \quad (7.37)$$

Simple estimate using  $s_j^2 = 1 + t_j^2$  shows that

$$|s_j t_k e^{i\sigma_k} - s_k t_j e^{i\sigma_j}|^2 \geq \frac{1}{4} (t_k - t_j)^2 \left( \frac{1}{1 + 2t_j^2} + \frac{1}{1 + 2t_k^2} \right) \geq \frac{1}{6} (t_k - t_j)^2. \quad (7.38)$$

Notice that the assumption  $\mathbf{t} \in \mathbb{I}^{W-1}$  was used only in the last inequality. By (7.37), (7.38) and the definition (7.1), Lemma 7.1 follows immediately.

**7.4. Proof of Lemma 7.2.** Let  $\mathbb{I}^c = \mathbb{R}_+ \setminus \mathbb{I}$ . Now we consider the domain sequence  $\vec{\mathbb{J}} = (\mathbb{J}_2, \dots, \mathbb{J}_W) \in \{\mathbb{I}, \mathbb{I}^c\}^{W-1}$ . We decompose the integral in Lemma 7.2 as follows

$$\int_{\mathbb{R}^{W-1} \setminus \mathbb{I}^{W-1}} \prod_{j=2}^W dt_j \mathbb{M}(\mathbf{t}) \mathbf{q}(\mathbf{t}) = \sum_{\substack{\vec{\mathbb{J}} \in \{\mathbb{I}, \mathbb{I}^c\}^{W-1} \\ \vec{\mathbb{J}} \neq \mathbb{I}^{W-1}}} \int_{\prod_{j=2}^W \mathbb{J}_j} \prod_{j=2}^W dt_j \mathbb{M}(\mathbf{t}) \mathbf{q}(\mathbf{t}). \quad (7.39)$$

Note the total number of the choices of such  $\vec{\mathbb{J}}$  in the sum above is  $2^{W-1} - 1$ . It suffices to consider one of these sequences  $\vec{\mathbb{J}} \in \{\mathbb{I}, \mathbb{I}^c\}^{W-1}$  in which there is at least one  $i$  such that  $\mathbb{J}_i = \mathbb{I}^c$ .

Recall the spanning tree  $\mathcal{G}_0 = (\mathcal{V}, \mathcal{E}_0)$  in Assumption 1.1. The simplest case is that there exists a linear spanning tree (a path)  $\mathcal{G}_0$  with

$$\mathcal{E}_0 = \{(i, i+1)\}_{i=1}^{W-1} \subset \mathcal{E}. \quad (7.40)$$

We first present the proof in this simplest case.

Now, we only keep the edges in the path  $\mathcal{E}_0$ , i.e. the terms with  $k = j - 1$  in (7.37), we also trivially discard the term  $1/(1 + 2t_j^2)$  from the sum  $1/(1 + 2t_{j-1}^2) + 1/(1 + 2t_j^2)$  in the estimate (7.38) (the first inequality), and finally we bound all  $M\mathfrak{A}(\hat{B})s_{j-1,j}/4$  by  $\mathfrak{L}$  defined in (7.3) from below. That means, we use the bound

$$\mathbb{M}(\mathbf{t}) \leq \prod_{j=2}^W \exp \left\{ -\mathfrak{L} \frac{(t_j - t_{j-1})^2}{1 + 2t_{j-1}^2} \right\} := \prod_{j=2}^W \check{\mathbb{M}}_j(\mathbf{t}). \quad (7.41)$$

Consequently, we have

$$\int_{\prod_{j=2}^W \mathbb{J}_j} \prod_{j=2}^W dt_j \mathbb{M}(\mathbf{t}) \mathbf{q}(\mathbf{t}) \leq \int_{\prod_{j=2}^W \mathbb{J}_j} \prod_{j=2}^W dt_j \prod_{j=2}^W t_j^{n_j} \check{\mathbb{M}}_j(\mathbf{t}). \quad (7.42)$$

Note that, as a function of  $\mathbf{t}$ ,  $\check{\mathbb{M}}_j(\mathbf{t})$  only depends on  $t_{j-1}$  and  $t_j$ .

Having fixed  $\vec{\mathbb{J}}$ , assume that  $k$  is the largest index such that  $\mathbb{J}_k = \mathbb{I}^c$ , i.e.  $t_{k+1}, \dots, t_W \in \mathbb{I}$ . Now, we claim that

$$\sum_{j=2}^W \frac{(t_j - t_{j-1})^2}{1 + 2t_{j-1}^2} \geq \sum_{j=2}^k \frac{(t_j - t_{j-1})^2}{1 + 2t_{j-1}^2} \geq \frac{1}{300k^2}, \quad \text{if } t_k \in \mathbb{I}^c. \quad (7.43)$$

To see (7.43), we use the following elementary facts

$$\frac{(t_j - t_{j-1})^2}{1 + 2t_{j-1}^2} \geq \frac{1}{3} \frac{(t_j - t_{j-1})^2}{t_{j-1}^2} = \frac{1}{3} (t_j/t_{j-1} - 1)^2, \quad \text{if } t_{j-1} \in \mathbb{I}^c \quad (7.44)$$

and

$$\frac{(t_j - t_{j-1})^2}{1 + 2t_{j-1}^2} \geq \frac{1}{3} (t_j - t_{j-1})^2, \quad \text{if } t_{j-1} \in \mathbb{I} \quad (7.45)$$

for all  $j = 2, \dots, W$ . We show (7.43) by contradiction. If (7.43) is violated, we have

$$\frac{(t_j - t_{j-1})^2}{1 + 2t_{j-1}^2} \leq \frac{1}{300k^2}, \quad \forall j = 2, \dots, k,$$

which together with (7.44) and (7.45) implies that

$$t_j \leq t_{j-1} \left(1 + \frac{1}{10k}\right) + \frac{1}{10k}. \quad (7.46)$$

Using (7.46) recursively yields

$$t_k \leq \left(1 + \frac{1}{10k}\right)^{k-1} (t_1 + 1) - 1 = \left(1 + \frac{1}{10k}\right)^{k-1} - 1 \leq \frac{1}{2}, \quad (7.47)$$

where in the second step we used the fact  $t_1 = 0$ . Note that (7.47) contradicts  $t_k \in \mathbb{I}^c$ . Hence, we verified (7.43).

Now, we split  $\prod_{j=2}^W \check{\mathbb{M}}_j(\mathbf{t})$  into two parts. We use one to control the integral, and the other will be estimated by (7.43). Specifically, substituting (7.43) into (7.42) we have

$$\int_{\prod_{j=2}^W \mathbb{J}_j} \prod_{j=2}^W dt_j \mathbb{M}(\mathbf{t}) \mathbf{q}(\mathbf{t}) \leq e^{-\frac{s}{600k^2}} \int_{\mathbb{R}_+^{W-1}} \prod_{j=2}^W dt_j \prod_{j=2}^W t_j^{n_j} (\check{\mathbb{M}}_j(\mathbf{t}))^{\frac{1}{2}}. \quad (7.48)$$

Therefore, what remains is to estimate the integral in (7.48), which can be done by elementary Gaussian integral step by step. More specifically, using (7.44) and (7.45) and the change of variable  $t_j/t_{j-1} - 1 \rightarrow t_j$  in case of  $t_{j-1} \in \mathbb{I}^c$  and  $t_j - t_{j-1} \rightarrow t_j$  in case of  $t_{j-1} \in \mathbb{I}$ , it is elementary to see that for any  $\ell = O(W)$ ,

$$\int_{\mathbb{R}_+} dt_j t_j^\ell (\check{\mathbb{M}}_j(\mathbf{t}))^{\frac{1}{2}} \leq \ell! \left(1 + c\mathfrak{L}^{-\frac{1}{2}}\right)^{O(\ell)} (t_{j-1}^{\ell+1} + 1) \leq e^{O(W \log N)} \left(1 + \mathfrak{L}^{-\frac{1}{2}}\right)^{O(\ell)} (t_{j-1}^{\ell+1} + 1). \quad (7.49)$$

Starting from  $j = W$ , using (7.49) to integrate (7.48) successively, the exponent of  $t_j$  increases linearly ( $n_j = O(1)$ ), thus we can get

$$\int_{\prod_{j=2}^W \mathbb{J}_j} \prod_{j=2}^W dt_j \mathbb{M}(\mathbf{t}) \mathbf{q}(\mathbf{t}) \leq e^{-\frac{s}{600W^2}} \cdot e^{O(W^2 \log N)} \cdot \left(1 + \mathfrak{L}^{-\frac{1}{2}}\right)^{O(W^2)}.$$

Then (7.4) follows from the definition of  $\mathfrak{L}$  in (7.3) and (5.38). Hence, we completed the proof for (7.4) when the spanning tree is given by (7.40).



Now, we consider more general spanning tree  $\mathcal{G}_0$  and regard 1 as its root. We start from the generalization of (7.41), namely,

$$\mathbb{M}(\mathbf{t}) \leq \prod_{\{i,j\} \in \mathcal{E}_0} \exp \left\{ -\mathfrak{L} \frac{(t_j - t_i)^2}{1 + 2t_i^2} \right\} := \prod_{\{i,j\} \in \mathcal{E}_0} \check{\mathbb{M}}_{i,j}(\mathbf{t}). \quad (7.50)$$

Here we make the convention that  $\text{dist}(1, i) = \text{dist}(1, j) - 1$  for all  $\{i, j\} \in \mathcal{E}_0$ , where  $\text{dist}(a, b)$  represents the distance between  $a$  and  $b$ . Now, if there is  $k'$  such that  $\mathbb{J}_{k'} \in \mathbb{I}^c$ , we can prove the following analogue of (7.43), namely,

$$\sum_{\{i,j\} \in \mathcal{E}_0} \frac{(t_j - t_i)^2}{1 + 2t_i^2} \geq \frac{1}{300k^2}$$

by performing the argument in (7.44)-(7.47) on the path connecting  $k'$  and the root 1. Consequently, we can get the analogue of (7.48) via replacing  $\check{\mathbb{M}}_j(\mathbf{t})$ 's by  $\check{\mathbb{M}}_{i,j}(\mathbf{t})$ 's. Finally, integrating  $t_j$ 's out successively, from the leaves to the root 1, yields the same conclusion, i.e. (7.4), for general  $\mathcal{G}_0$ . Therefore, we completed the proof of Lemma 7.2.

## 8. GAUSSIAN MEASURE IN THE VICINITIES

From now on, we can restrict ourselves to the Type I, II and III vicinities. As a preparation of the proofs of Lemmas 5.8 and 5.9, we will show in this section that the exponential function

$$\exp \left\{ -M(\hat{K}(\hat{X}, V) + \hat{L}(\hat{B}, T)) \right\} \quad (8.1)$$

is approximately a Gaussian measure (unnormalized).

**8.1. Parametrization and initial approximation in the vicinities.** We change the  $\mathbf{x}$ ,  $\mathbf{b}$ ,  $\mathbf{t}$ ,  $\mathbf{v}$ -variables to a new set of variables, namely,  $\mathring{\mathbf{x}}$ ,  $\mathring{\mathbf{b}}$ ,  $\mathring{\mathbf{t}}$  and  $\mathring{\mathbf{v}}$ . The precise definition of  $\mathring{x}$  differs in the different vicinities. To distinguish the parameterization, we set  $\varkappa = \pm, +$ , or  $-$ , corresponding to Type I, II or III vicinity, respectively. Recalling  $D_\varkappa$  from (1.24). For each  $j$  and each  $\varkappa$ , we then set

$$\begin{aligned} \hat{X}_j &= D_\varkappa \text{diag} \left( \exp \{ \mathbf{i} \mathring{x}_{j,1} / \sqrt{M} \}, \exp \{ \mathbf{i} \mathring{x}_{j,2} / \sqrt{M} \} \right), \quad \mathring{x}_{j,a} / \sqrt{M} \in [-\pi, \pi], \\ \hat{B}_j &= D_\pm + D_\pm \text{diag} \left( \mathring{b}_{j,1} / \sqrt{M}, \mathring{b}_{j,2} / \sqrt{M} \right), \quad t_j = \mathring{t}_j / \sqrt{M}. \end{aligned} \quad (8.2)$$

If  $\varkappa = \pm$ , we also need to parameterize  $v_j$  by

$$v_j = \mathring{v}_j / \sqrt{M}. \quad (8.3)$$

We set the vectors

$$\begin{aligned} \mathring{\mathbf{b}}_a &:= (\mathring{b}_{1,a}, \dots, \mathring{b}_{W,a}), \quad \mathring{\mathbf{x}}_a := (\mathring{x}_{1,a}, \dots, \mathring{x}_{W,a}), \quad a = 1, 2, \\ \mathring{\mathbf{t}} &:= (\mathring{t}_2, \dots, \mathring{t}_W), \quad \mathring{\mathbf{v}} := (\mathring{v}_2, \dots, \mathring{v}_W). \end{aligned}$$

Accordingly, recalling the quantity  $\Theta$  from (5.33), we introduce the domains

$$\mathring{\Upsilon} \equiv \mathring{\Upsilon}(N, \varepsilon_0) := \{\mathbf{a} \in \mathbb{R}^W : \|\mathbf{a}\|_2^2 \leq \Theta\}, \quad \mathring{\Upsilon}_S \equiv \mathring{\Upsilon}_S(N, \varepsilon_0) := \{\mathbf{a} \in \mathbb{R}_+^{W-1} : -\mathbf{a}' S^{(1)} \mathbf{a} \leq \Theta\}.$$

We remind here, as mentioned above, in the sequel, the small constant  $\varepsilon_0$  in  $\mathring{\Upsilon}$  and  $\mathring{\Upsilon}_S$  may be different from line to line, subject to (5.37). Now, by the definition of the Type I', II and III vicinities in Definition 5.5 and the parametrization in (8.2) and (8.3), we can redefine the vicinities as follows.

**Definition 8.1.** *We can redefine three types of vicinities as follows.*

- *Type I' vicinity :*  $(\mathring{\mathbf{b}}_1, \mathring{\mathbf{b}}_2, \mathring{\mathbf{x}}_1, \mathring{\mathbf{x}}_2, \mathring{\mathbf{t}}, \mathring{\mathbf{v}}) \in \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon}_S \times \mathring{\Upsilon}_S$ , with  $\varkappa = \pm$ .
- *Type II vicinity :*  $(\mathring{\mathbf{b}}_1, \mathring{\mathbf{b}}_2, \mathring{\mathbf{x}}_1, \mathring{\mathbf{x}}_2, \mathring{\mathbf{t}}, \mathbf{v}) \in \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon}_S \times \mathbb{I}^{W-1}$ , with  $\varkappa = +$ .
- *Type III vicinity :*  $(\mathring{\mathbf{b}}_1, \mathring{\mathbf{b}}_2, \mathring{\mathbf{x}}_1, \mathring{\mathbf{x}}_2, \mathring{\mathbf{t}}, \mathbf{v}) \in \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon} \times \mathring{\Upsilon}_S \times \mathbb{I}^{W-1}$ , with  $\varkappa = -$ .

We recall from (7.6) the fact

$$\mathbf{t} \in \hat{\Upsilon}_S \implies \|\mathbf{t}\|_\infty = O(\Theta). \quad (8.4)$$

Now, we use the representation (5.2). Then, for the Type I vicinity, we change  $\mathbf{x}, \mathbf{b}, \mathbf{t}, \mathbf{v}$ -variables to  $\mathring{\mathbf{x}}, \mathring{\mathbf{b}}, \mathring{\mathbf{t}}, \mathring{\mathbf{v}}$ -variables according to (8.2) with  $\varkappa = \pm$ , thus

$$\begin{aligned} 2^W \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \Upsilon_S) &= \frac{M^2}{(n!)^2 4^W \pi^{2W+4}} \int_{\mathbb{L}^{2W-2}} \prod_{j=2}^W \frac{d\theta_j}{2\pi} \prod_{j=2}^W \frac{d\sigma_j}{2\pi} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\mathring{b}_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\mathring{b}_{j,2} \\ &\times \int_{\hat{\Upsilon}} \prod_{j=1}^W d\mathring{x}_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\mathring{x}_{j,2} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W 2\mathring{t}_j d\mathring{t}_j \int_{\hat{\Upsilon}_S} \prod_{j=2}^W 2\mathring{v}_j d\mathring{v}_j \prod_{j=1}^W \exp\left\{i \frac{\mathring{x}_{j,1} + \mathring{x}_{j,2}}{\sqrt{M}}\right\} \\ &\times \exp\left\{-M(\mathring{K}(\hat{X}, V) + \mathring{L}(\hat{B}, T))\right\} \cdot \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot \mathbf{A}(\hat{X}, \hat{B}, V, T). \end{aligned} \quad (8.5)$$

For the Type II or III vicinities, i.e.  $\varkappa = +$  or  $-$ , we change  $\mathbf{x}, \mathbf{b}, \mathbf{t}$ -variables to  $\mathring{\mathbf{x}}, \mathring{\mathbf{b}}, \mathring{\mathbf{t}}$ -variables. Consequently, we have

$$\begin{aligned} \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_\varkappa^x, \Upsilon_\varkappa^x, \Upsilon_S, \mathbb{L}^{W-1}) &= \frac{(-a_\varkappa^2)^W}{(n!)^2} \cdot \frac{M^{W+1}}{8^W \pi^{2W+4}} \cdot \int_{\mathbb{L}^{2W-2}} \prod_{j=2}^W \frac{d\theta_j}{2\pi} \prod_{j=2}^W \frac{d\sigma_j}{2\pi} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\mathring{b}_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\mathring{b}_{j,2} \\ &\times \int_{\hat{\Upsilon}} \prod_{j=1}^W d\mathring{x}_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\mathring{x}_{j,2} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W 2\mathring{t}_j d\mathring{t}_j \int_{\mathbb{L}^{W-1}} \prod_{j=2}^W 2v_j dv_j \prod_{j=1}^W \exp\left\{i \frac{\mathring{x}_{j,1} + \mathring{x}_{j,2}}{\sqrt{M}}\right\} \\ &\times \exp\left\{-M(\mathring{K}(\hat{X}, V) + \mathring{L}(\hat{B}, T))\right\} \cdot \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot \mathbf{A}(\hat{X}, \hat{B}, V, T). \end{aligned} \quad (8.6)$$

We will also need the following facts

$$\prod_{j=1}^W |(x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2| = e^{O(W)}, \quad |\mathbf{A}(\hat{X}, \hat{B}, V, T)| \leq e^{O(WN^{\varepsilon_2})} \quad (8.7)$$

if

$$\mathbf{x}_1, \mathbf{x}_2 \in \hat{\Sigma}^W, \quad b_{j,1} = a_+ + o(1), \quad b_{j,2} = -a_- + o(1), \quad t_j = o(1), \quad \forall j = 1, \dots, N,$$

which always hold in these types of vicinities. The first estimate in (8.7) is trivial, and the second follows from Lemma 6.1.

Now, we approximate (8.1) in the vicinities. For any  $\vartheta \in \mathbb{L}$ , we introduce the matrices

$$\mathcal{E}_+(\vartheta) := \begin{pmatrix} 0 & e^{i\vartheta} \\ e^{-i\vartheta} & 0 \end{pmatrix}, \quad \mathcal{E}_-(\vartheta) := \begin{pmatrix} 0 & e^{i\vartheta} \\ -e^{-i\vartheta} & 0 \end{pmatrix}.$$

Then, with the parameterization above, expanding  $\hat{X}_j$  in (3.22) and  $T_j$  in (3.24) up to the second order, we can write

$$\begin{aligned} \hat{X}_j &= D_\varkappa + \frac{\mathbf{i}}{\sqrt{M}} D_\varkappa \text{diag}(\mathring{x}_{j,1}, \mathring{x}_{j,2}) + \frac{1}{M} R_j^x, \quad \varkappa = \pm, +, -, \\ T_j &= I + \frac{\mathring{t}_j}{\sqrt{M}} \mathcal{E}_+(\sigma_j) + \frac{1}{M} R_j^t. \end{aligned} \quad (8.8)$$

For  $\varkappa = \pm$ , we also expand  $V_j$  in (3.24) up to the second order, namely,

$$V_j = I + \frac{\mathring{v}_j}{\sqrt{M}} \mathcal{E}_-(\theta_j) + \frac{1}{M} R_j^v. \quad (8.9)$$

We just take (8.8) and (8.9) as the definition of  $R_j^x$ ,  $R_j^t$  and  $R_j^v$ . Note that  $R_j^x$  is actually  $\varkappa$ -dependent. However, this dependence is irrelevant for our analysis thus is suppressed from the notation. It is elementary that

$$\|R_j^x\|_{\max} = O(\dot{x}_{j,1}^2 + \dot{x}_{j,2}^2), \quad \|R_j^t\|_{\max} = O(\dot{t}_j^2), \quad \|R_j^v\|_{\max} = O(\dot{v}_j^2). \quad (8.10)$$

Here  $\|\cdot\|_{\max}$  represents the max-norm of a matrix.

Recall the facts (5.10) and (5.19)

$$\begin{aligned} M\dot{L}(\hat{B}, T) &= M\left(\dot{\ell}_{++}(\mathbf{b}_1) + \dot{\ell}_{--}(\mathbf{b}_2)\right) + M\ell_S(\hat{B}, T), \\ M\dot{K}(\hat{X}, V) &= M\left(-\dot{\ell}_{++}(\mathbf{x}_1) - \dot{\ell}_{+-}(\mathbf{x}_2)\right) + M\ell_S(\hat{X}, V). \end{aligned} \quad (8.11)$$

In light of (5.17)-(5.19), we can also represent  $M\dot{K}(\hat{X}, V)$  in the following two alternative ways

$$M\dot{K}(\hat{X}, V) = M\left(-\dot{\ell}_{++}(\mathbf{x}_1) - \dot{\ell}_{++}(\mathbf{x}_2)\right) + M\ell_S(\hat{X}, V) + M(K(D_+, I) - K(D_{\pm}, I)), \quad (8.12)$$

$$M\dot{K}(\hat{X}, V) = M\left(-\dot{\ell}_{+-}(\mathbf{x}_1) - \dot{\ell}_{+-}(\mathbf{x}_2)\right) + M\ell_S(\hat{X}, V) + M(K(D_-, I) - K(D_{\pm}, I)). \quad (8.13)$$

We will use three representations of  $M\dot{K}(\hat{X}, V)$  in (8.11), (8.12) and (8.13) for Type I', II and III vicinities respectively. In addition, we introduce the matrices

$$\mathbb{A}_+ := (1 + a_+^2)I + a_+^2 S, \quad \mathbb{A}_- := (1 + a_-^2)I + a_-^2 S. \quad (8.14)$$

Then, we have the following lemma.

**Lemma 8.2.** *With the parametrization in (8.8), we have the following approximations.*

- Let  $\mathring{\mathbf{b}}_1, \mathring{\mathbf{b}}_2 \in \mathbb{C}^W$  and  $\|\mathring{\mathbf{b}}_1\|_{\infty}, \|\mathring{\mathbf{b}}_2\|_{\infty} = o(\sqrt{M})$ , we have

$$M\left(\dot{\ell}_{++}(\mathbf{b}_1) + \dot{\ell}_{--}(\mathbf{b}_2)\right) = \frac{1}{2}\mathring{\mathbf{b}}_1' \mathbb{A}_+ \mathring{\mathbf{b}}_1 + \frac{1}{2}\mathring{\mathbf{b}}_2' \mathbb{A}_- \mathring{\mathbf{b}}_2 + R^b, \quad R^b = O\left(\frac{\sum_{a=1,2} \|\mathring{\mathbf{b}}_a\|_3^3}{\sqrt{M}}\right). \quad (8.15)$$

- Let  $\varkappa = \pm$  and  $\mathring{\mathbf{x}}_1, \mathring{\mathbf{x}}_2 \in \mathbb{C}^W$  and  $\|\mathring{\mathbf{x}}_1\|_{\infty}, \|\mathring{\mathbf{x}}_2\|_{\infty} = o(\sqrt{M})$ , we have

$$M\left(-\dot{\ell}_{++}(\mathbf{x}_1) - \dot{\ell}_{+-}(\mathbf{x}_2)\right) = \frac{1}{2}\mathring{\mathbf{x}}_1' \mathbb{A}_+ \mathring{\mathbf{x}}_1 + \frac{1}{2}\mathring{\mathbf{x}}_2' \mathbb{A}_- \mathring{\mathbf{x}}_2 + R_{\pm}^x, \quad R_{\pm}^x = O\left(\frac{\sum_{a=1,2} \|\mathring{\mathbf{x}}_a\|_3^3}{\sqrt{M}}\right). \quad (8.16)$$

- In the Type II vicinity, we have

$$M\left(-\dot{\ell}_{++}(\mathbf{x}_1) - \dot{\ell}_{++}(\mathbf{x}_2)\right) = \frac{1}{2}\mathring{\mathbf{x}}_1' \mathbb{A}_+ \mathring{\mathbf{x}}_1 + \frac{1}{2}\mathring{\mathbf{x}}_2' \mathbb{A}_+ \mathring{\mathbf{x}}_2 + R_+^x, \quad R_+^x = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right). \quad (8.17)$$

- In the Type III vicinity, we have

$$M\left(-\dot{\ell}_{+-}(\mathbf{x}_1) - \dot{\ell}_{+-}(\mathbf{x}_2)\right) = \frac{1}{2}\mathring{\mathbf{x}}_1' \mathbb{A}_- \mathring{\mathbf{x}}_1 + \frac{1}{2}\mathring{\mathbf{x}}_2' \mathbb{A}_- \mathring{\mathbf{x}}_2 + R_-^x, \quad R_-^x = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right). \quad (8.18)$$

Here  $R^b$ ,  $R_{\pm}^x$ ,  $R_+^x$  and  $R_-^x$  are remainder terms of the Taylor expansion of the function  $\ell(\mathbf{a})$  defined in (5.6).

**Remark 8.3.** Here we stated (8.15) and (8.16) in the domains much larger than the Type I' vicinity for further discussion. In addition, the restriction  $\|\mathring{\mathbf{b}}_a\|_{\infty}$  and  $\|\mathring{\mathbf{x}}_a\|_{\infty}$  for  $a = 1, 2$  is imposed to avoid the ambiguity of the definition of the logarithmic term in the function  $\ell(\mathbf{a})$ .

*Proof.* It follows from the Taylor expansion of the function  $\ell(\mathbf{a})$  easily.  $\square$

Then, according to (8.11)-(8.13), what remains is to approximate  $M\ell_S(\hat{B}, T)$  and  $M\ell_S(\hat{X}, V)$  in the vicinities. Recalling the definition in (5.6) and the parameterization in (8.2), we can rewrite

$$\begin{aligned} M\ell_S(\hat{B}, T) &= \frac{1}{2} \sum_{j,k} \mathfrak{s}_{jk} |s_j \mathring{t}_k e^{i\sigma_k} - s_k \mathring{t}_j e^{i\sigma_j}|^2 \\ &\quad \times \left( a_+ - a_- + \frac{a_+ \mathring{b}_{j,1} - a_- \mathring{b}_{j,2}}{\sqrt{M}} \right) \left( a_+ - a_- + \frac{a_+ \mathring{b}_{k,1} - a_- \mathring{b}_{k,2}}{\sqrt{M}} \right) \\ &=: \frac{(a_+ - a_-)^2}{2} \sum_{j,k} \mathfrak{s}_{jk} |\mathring{t}_k e^{i\sigma_k} - \mathring{t}_j e^{i\sigma_j}|^2 + R^{t,b}. \end{aligned} \quad (8.19)$$

We take the above equation as the definition of  $R^{t,b}$ . Now, we set

$$\tau_{j,1} := \mathring{t}_j \cos \sigma_j, \quad \tau_{j,2} := \mathring{t}_j \sin \sigma_j, \quad \forall j = 2, \dots, W$$

and change the variables and the measure as

$$(\mathring{t}_j, \sigma_j) \rightarrow (\tau_{j,1}, \tau_{j,2}), \quad 2\mathring{t}_j d\mathring{t}_j \frac{d\sigma_j}{2\pi} \rightarrow \frac{1}{\pi} d\tau_{j,1} d\tau_{j,2}. \quad (8.20)$$

In the Type I' vicinity, we can do the same thing for  $M\ell_S(\hat{X}, V)$ , namely,

$$M\ell_S(\hat{X}, V) =: \frac{(a_+ - a_-)^2}{2} \sum_{j,k} \mathfrak{s}_{jk} |\mathring{v}_k e^{i\theta_k} - \mathring{v}_j e^{i\theta_j}|^2 + R_{\pm}^{v,x}, \quad (8.21)$$

where  $R_{\pm}^{v,x}$  is the remainder term. Then we set

$$v_{j,1} := \mathring{v}_j \cos \theta_j, \quad v_{j,2} := \mathring{v}_j \sin \theta_j, \quad \forall j = 2, \dots, W$$

and change the variables and measure as

$$(\mathring{v}_j, \theta_j) \rightarrow (v_{j,1}, v_{j,2}), \quad 2\mathring{v}_j d\mathring{v}_j \frac{d\theta_j}{2\pi} \rightarrow \frac{1}{\pi} dv_{j,1} dv_{j,2}. \quad (8.22)$$

Now, we introduce the vectors

$$\boldsymbol{\tau}_a = (\tau_{2,a}, \dots, \tau_{W,a}), \quad \boldsymbol{v}_a = (v_{2,a}, \dots, v_{W,a}), \quad a = 1, 2.$$

With this notation, we can rewrite (8.19) and (8.21) as

$$\begin{aligned} M\ell_S(\hat{B}, T) &= -(a_+ - a_-)^2 \sum_{a=1,2} \boldsymbol{\tau}'_a S^{(1)} \boldsymbol{\tau}_a + R^{t,b}, \\ M\ell_S(\hat{X}, V) &= -(a_+ - a_-)^2 \sum_{a=1,2} \boldsymbol{v}'_a S^{(1)} \boldsymbol{v}_a + R_{\pm}^{v,x}. \end{aligned} \quad (8.23)$$

According to (8.20) and (8.22), we can express (8.5) as an integral over  $\mathring{\mathbf{b}}$ ,  $\mathring{\mathbf{x}}$ ,  $\mathring{\boldsymbol{\tau}}$  and  $\mathring{\mathbf{v}}$ -variables. However, we need to specify the domains of  $\mathring{\boldsymbol{\tau}}$  and  $\mathring{\mathbf{v}}$ -variables in advance. Our aim is to restrict the integral in the domains

$$\boldsymbol{\tau}_a \in \mathring{\Upsilon}_S, \quad \boldsymbol{v}_a \in \mathring{\Upsilon}_S, \quad a = 1, 2. \quad (8.24)$$

Taking  $\mathring{\mathbf{t}}$  for instance, we see that

$$(t_j - t_k)^2 \leq |t_j e^{i\sigma_j} - t_k e^{i\sigma_k}|^2 = (\tau_{j,1} - \tau_{k,1})^2 + (\tau_{j,2} - \tau_{k,2})^2,$$

which actually implies

$$\boldsymbol{\tau}_a \in \mathring{\Upsilon}_S \quad \text{for } a = 1, 2 \implies \mathring{\mathbf{t}} \in \mathring{\Upsilon}_S. \quad (8.25)$$

However the reverse of (8.25) may not be true. That means, (8.24) is stronger than  $(\mathring{\mathbf{t}}, \mathring{\mathbf{v}}, \boldsymbol{\sigma}, \boldsymbol{\theta}) \in \mathring{\Upsilon}_S \times \mathring{\Upsilon}_S \times \mathbb{L}^{W-1} \times \mathbb{L}^{W-1}$ . To show the truncation to (8.24) from  $(\mathring{\mathbf{t}}, \mathring{\mathbf{v}}, \boldsymbol{\sigma}, \boldsymbol{\theta}) \in \mathring{\Upsilon}_S \times \mathring{\Upsilon}_S \times \mathbb{L}^{W-1} \times \mathbb{L}^{W-1}$

is harmless in the integral (8.5), we need to bound  $R^{t,b}$  and  $R_{\pm}^{v,x}$  in terms of  $\tau'_a S^{(1)} \tau_a$  and  $\mathbf{v}'_a S^{(1)} \mathbf{v}_a$ , respectively. More specifically, we need the following lemma.

**Lemma 8.4.** *In the Type I' vicinity, we have*

$$\begin{aligned} |R^{t,b}| &\leq O\left(\frac{\Theta^{\frac{1}{2}}}{\sqrt{M}}\right) \sum_{a=1,2} (-\tau'_a S^{(1)} \tau_a) + O\left(\frac{\Theta^{\frac{7}{2}}}{M}\right) \sum_{a=1,2} (-\tau'_a S^{(1)} \tau_a)^{1/2} + O\left(\frac{\Theta^7}{M^2}\right), \\ |R_{\pm}^{v,x}| &\leq O\left(\frac{\Theta^{\frac{1}{2}}}{\sqrt{M}}\right) \sum_{a=1,2} (-\mathbf{v}'_a S^{(1)} \mathbf{v}_a) + O\left(\frac{\Theta^{\frac{7}{2}}}{M}\right) \sum_{a=1,2} (-\mathbf{v}'_a S^{(1)} \mathbf{v}_a)^{1/2} + O\left(\frac{\Theta^7}{M^2}\right). \end{aligned} \quad (8.26)$$

*Proof.* Since the proofs of these two bounds are nearly the same, we only state the details for the first one. By (8.4), we have

$$s_j = 1 + O\left(\frac{\Theta^2}{M}\right), \quad \forall j = 2, \dots, W.$$

Then it is not difficult to see that

$$\begin{aligned} |s_j \dot{t}_k e^{i\sigma_k} - s_k \dot{t}_j e^{i\sigma_j}|^2 &= \left| \dot{t}_k e^{i\sigma_k} - \dot{t}_j e^{i\sigma_j} + O\left(\frac{\Theta^3}{M}\right) \right|^2 \\ &= \sum_{a=1,2} (\tau_{j,a} - \tau_{k,a})^2 + O\left(\frac{\Theta^3}{M}\right) \sum_{a=1,2} |\tau_{j,a} - \tau_{k,a}| + O\left(\frac{\Theta^6}{M^2}\right). \end{aligned} \quad (8.27)$$

Now, by the fact from Definition 8.1

$$\|\mathring{\mathbf{b}}_a\|_{\infty} \leq \|\mathring{\mathbf{b}}_a\|_2 = O(\Theta^{\frac{1}{2}}) \quad \text{for } a = 1, 2, \quad (8.28)$$

we have

$$\left(a_+ - a_- + \frac{a_+ \dot{b}_{j,1} - a_- \dot{b}_{j,2}}{\sqrt{M}}\right) \left(a_+ - a_- + \frac{a_+ \dot{b}_{k,1} - a_- \dot{b}_{k,2}}{\sqrt{M}}\right) = (a_+ - a_-)^2 + O\left(\frac{\Theta^{\frac{1}{2}}}{\sqrt{M}}\right). \quad (8.29)$$

Substituting (8.27) and (8.29) into (8.19) yields

$$\begin{aligned} |R^{t,b}| &\leq O\left(\frac{\Theta^{\frac{1}{2}}}{\sqrt{M}}\right) \sum_{a=1,2} \sum_{j,k} \mathfrak{s}_{jk} (\tau_{j,a} - \tau_{k,a})^2 + O\left(\frac{\Theta^3}{M}\right) \sum_{a=1,2} \sum_{j,k} \mathfrak{s}_{jk} |\tau_{j,a} - \tau_{k,a}| + O\left(\frac{\Theta^7}{M^2}\right) \\ &\leq O\left(\frac{\Theta^{\frac{1}{2}}}{\sqrt{M}}\right) \sum_{a=1,2} (-\tau'_a S^{(1)} \tau_a) + O\left(\frac{\Theta^{\frac{7}{2}}}{M}\right) \sum_{a=1,2} (-\tau'_a S^{(1)} \tau_a)^{1/2} + O\left(\frac{\Theta^7}{M^2}\right). \end{aligned}$$

Here we used Cauchy Schwarz inequality and  $\sum_{j,k} \mathfrak{s}_{jk} = O(W) = O(\Theta)$  in the second step. The bound for  $|R_{\pm}^{v,x}|$  can be proved analogously. Hence, we completed the proof of Lemma 8.4.  $\square$

Roughly speaking, by (8.23) and Lemma 8.4 we have

$$\begin{aligned} M\ell_S(\hat{B}, T) &= -(a_+ - a_- + o(1))^2 \sum_{a=1,2} \tau'_a S^{(1)} \tau_a + o\left(\sum_{a=1,2} (-\tau'_a S^{(1)} \tau_a)^{1/2}\right) + o(1), \\ M\ell_S(\hat{X}, V) &= -(a_+ - a_- + o(1))^2 \sum_{a=1,2} \mathbf{v}'_a S^{(1)} \mathbf{v}_a + o\left(\sum_{a=1,2} (-\mathbf{v}'_a S^{(1)} \mathbf{v}_a)^{1/2}\right) + o(1). \end{aligned}$$

Then it is obvious that if one of  $\tau_1, \tau_2, \mathbf{v}_1$  and  $\mathbf{v}_2$  is not in  $\mathring{\Upsilon}_S$ , we will get (7.26). Hence, using (8.7), we can discard the integral outside the vicinity, analogously to the proof of Lemma 7.7. More specifically, in the sequel, we can and do assume

$$\tau_1, \tau_2, \mathbf{v}_1, \mathbf{v}_2 \in \mathring{\Upsilon}_S. \quad (8.30)$$

Now, plugging (8.30) into (8.26) in turn yields the bound

$$|R^{t,b}| = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}} \vee \frac{\Theta^4}{M}\right), \quad |R_{\pm}^{v,x}| = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}} \vee \frac{\Theta^4}{M}\right). \quad (8.31)$$

By the discussion above, for the Type I vicinity, we can write (8.5) as

$$\begin{aligned}
2^W \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \Upsilon_S) &= \frac{M^2}{(n!)^2 4^W \pi^{4W+2}} \cdot \int_{\dot{\Upsilon}} \prod_{j=1}^W d\dot{b}_{j,1} \int_{\dot{\Upsilon}} \prod_{j=1}^W d\dot{b}_{j,2} \\
&\times \int_{\dot{\Upsilon}} \prod_{j=1}^W d\dot{x}_{j,1} \int_{\dot{\Upsilon}} \prod_{j=1}^W d\dot{x}_{j,2} \int_{\dot{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,1} \int_{\dot{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,2} \int_{\dot{\Upsilon}_S} \prod_{j=2}^W dv_{j,1} \int_{\dot{\Upsilon}_S} \prod_{j=2}^W dv_{j,2} \\
&\times \exp \left\{ -\frac{1}{2} \dot{\mathbf{b}}_1' \mathbb{A}_+ \dot{\mathbf{b}}_1 - \frac{1}{2} \dot{\mathbf{b}}_2' \mathbb{A}_- \dot{\mathbf{b}}_2 - R^b \right\} \cdot \exp \left\{ -\frac{1}{2} \dot{\mathbf{x}}_1' \mathbb{A}_+ \dot{\mathbf{x}}_1 - \frac{1}{2} \dot{\mathbf{x}}_2' \mathbb{A}_- \dot{\mathbf{x}}_2 - R_{\pm}^x \right\} \\
&\times \exp \left\{ (a_+ - a_-)^2 \sum_{a=1,2} \tau_a' S^{(1)} \tau_a - R^{t,b} \right\} \cdot \exp \left\{ (a_+ - a_-)^2 \sum_{a=1,2} \mathbf{v}_a' S^{(1)} \mathbf{v}_a - R_{\pm}^{v,x} \right\} \\
&\times \prod_{j=1}^W \exp \left\{ i \frac{\dot{x}_{j,1} + \dot{x}_{j,2}}{\sqrt{M}} \right\} \cdot \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot \mathbf{A}(\hat{X}, \hat{B}, V, T) + O(e^{-\Theta}), \quad (8.32)
\end{aligned}$$

where the error term stems from the truncation of the vicinity  $(\dot{\mathbf{t}}, \dot{\mathbf{v}}, \boldsymbol{\sigma}, \boldsymbol{\theta}) \in \dot{\Upsilon}_S \times \dot{\Upsilon}_S \times \mathbb{L}^{W-1} \times \mathbb{L}^{W-1}$  to  $(\boldsymbol{\tau}_1, \boldsymbol{\tau}_2, \mathbf{v}_1, \mathbf{v}_2) \in \dot{\Upsilon}_S \times \dot{\Upsilon}_S \times \dot{\Upsilon}_S \times \dot{\Upsilon}_S$ .

Now, for the Type II and III vicinities, the discussion on  $\ell_S(\hat{B}, T)$  is of course the same. For  $\ell_S(\hat{X}, V)$ , we make the following approximation. For the Type II vicinity, using the notation in (5.21), we can write

$$\begin{aligned}
M \ell_S(\hat{X}, V) &= \frac{-M a_+^2}{2} \sum_{j,k} \mathbf{s}_{jk}^v \left( \frac{\dot{x}_{j,1}}{\sqrt{M}} - \frac{\dot{x}_{j,2}}{\sqrt{M}} + O\left(\frac{\dot{x}_{j,1}^2 + \dot{x}_{j,2}^2}{M}\right) \right) \left( \frac{\dot{x}_{k,1}}{\sqrt{M}} - \frac{\dot{x}_{k,2}}{\sqrt{M}} + O\left(\frac{\dot{x}_{k,1}^2 + \dot{x}_{k,2}^2}{M}\right) \right) \\
&=: -\frac{a_+^2}{2} \sum_{j,k} \mathbf{s}_{jk}^v (\dot{x}_{j,1} - \dot{x}_{j,2})(\dot{x}_{k,1} - \dot{x}_{k,2}) + R_+^{v,x}. \quad (8.33)
\end{aligned}$$

It is easy to see that

$$R_+^{v,x} = O\left(\frac{\|\dot{\mathbf{x}}_1\|_3^3 + \|\dot{\mathbf{x}}_2\|_3^3}{\sqrt{M}}\right) = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right), \quad \text{for } \dot{\mathbf{x}}_1, \dot{\mathbf{x}}_2 \in \dot{\Upsilon}. \quad (8.34)$$

Combining (5.15), (5.18), (8.17) and (8.33) we obtain

$$\begin{aligned}
M(K(\hat{X}, V) - K(D_+, I)) &= \frac{1}{2} \dot{\mathbf{x}}_1' \mathbb{A}_+ \dot{\mathbf{x}}_1 + \frac{1}{2} \dot{\mathbf{x}}_2' \mathbb{A}_- \dot{\mathbf{x}}_2 - \frac{a_+^2}{2} \sum_{j,k} \mathbf{s}_{jk}^v (\dot{x}_{j,1} - \dot{x}_{j,2})(\dot{x}_{k,1} - \dot{x}_{k,2}) + R_+^x + R_+^{v,x} \\
&= \frac{1}{2} \dot{\mathbf{x}}' \mathbb{A}_+^v \dot{\mathbf{x}} + R_+^x + R_+^{v,x},
\end{aligned}$$

where

$$\dot{\mathbf{x}} := (\dot{\mathbf{x}}_1', \dot{\mathbf{x}}_2')', \quad \mathbb{A}_+^v := (1 + a_+^2) I_{2W} + a_+^2 \mathbb{S}^v, \quad (8.35)$$

and recall that  $\mathbb{S}^v$  is defined in (5.22). Analogously, for the Type III vicinity, we can write

$$M(K(\hat{X}, V) - K(D_-, I)) = \frac{1}{2} \dot{\mathbf{x}}' \mathbb{A}_-^v \dot{\mathbf{x}} + R_-^x + R_-^{v,x}, \quad (8.36)$$

where

$$\mathbb{A}_-^v := (1 + a_-^2) I_{2W} + a_-^2 \mathbb{S}^v, \quad R_-^{v,x} = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right). \quad (8.37)$$

Consequently, by (8.12) and (8.13) we can write (8.6) for  $\kappa = +, -$  as

$$\begin{aligned}
\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_{\mathcal{Z}}^x, \Upsilon_{\mathcal{Z}}^x, \Upsilon_S, \mathbb{L}^{W-1}) &= \exp \left\{ M(K(D_{\pm}, I) - K(D_{\mathcal{Z}}, I)) \right\} \cdot \frac{(-a_{\mathcal{Z}}^2)^W}{(n!)^2} \cdot \frac{M^{W+1}}{8^W \pi^{3W+3}} \\
&\times \int_{\mathbb{L}^{W-1}} \prod_{j=2}^W \frac{d\theta_j}{2\pi} \int_{\mathbb{L}^{W-1}} \prod_{j=2}^W 2v_j dv_j \int_{\dot{\Upsilon}} \prod_{j=1}^W d\dot{b}_{j,1} \int_{\dot{\Upsilon}} \prod_{j=1}^W d\dot{b}_{j,2} \int_{\dot{\Upsilon}} \prod_{j=1}^W d\dot{x}_{j,1} \int_{\dot{\Upsilon}} \prod_{j=1}^W d\dot{x}_{j,2}
\end{aligned}$$

$$\begin{aligned}
& \times \int_{\hat{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,1} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,2} \cdot \exp \left\{ -\frac{1}{2} \mathring{\mathbf{b}}'_1 \mathbb{A}_+ \mathring{\mathbf{b}}_1 - \frac{1}{2} \mathring{\mathbf{b}}'_2 \mathbb{A}_- \mathring{\mathbf{b}}_2 - R^b \right\} \\
& \times \exp \left\{ -\frac{1}{2} \mathring{\mathbf{x}}' \mathbb{A}_{\mathcal{K}}^v \mathring{\mathbf{x}} - R_{\mathcal{K}}^x - R_{\mathcal{K}}^{v,x} \right\} \cdot \exp \left\{ (a_+ - a_-)^2 \sum_{a=1,2} \tau'_a S^{(1)} \tau_a - R^{t,b} \right\} \\
& \times \prod_{j=1}^W \exp \left\{ i \frac{\mathring{x}_{j,1} + \mathring{x}_{j,2}}{\sqrt{M}} \right\} \cdot \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot A(\hat{X}, \hat{B}, V, T) + O(e^{-\Theta}). \quad (8.38)
\end{aligned}$$

**8.2. Steepest descent paths in the vicinities.** In order to estimate the integrals (8.32) and (8.38) properly, we need to control various remainder terms in (8.32) and (8.38) to reduce these integrals to Gaussian ones. The final result is collected in Proposition 8.6 at the end of this section. As a preparation, we shall further deform the contours of  $\mathring{\mathbf{b}}$ -variables and  $\mathring{\mathbf{x}}$ -variables to the steepest descent paths. We mainly provide the discussion for the  $\mathring{\mathbf{b}}$ -variables, that for the  $\mathring{\mathbf{x}}$ -variables is analogous.

For simplicity, in this section, we assume  $0 \leq E \leq \sqrt{2} - \kappa$ , the case  $-\sqrt{2} + \kappa \leq E \leq 0$  can be discussed similarly. We introduce the eigendecomposition of  $S$  as

$$S = \mathbf{U} \hat{S} \mathbf{U}'.$$

Note that  $\mathbf{U}$  is an orthogonal matrix thus the entries are all real. Now, we perform the change of coordinate

$$\mathbf{c}_a = (c_{1,a}, \dots, c_{W,a})' := \mathbf{U}' \mathring{\mathbf{b}}_a, \quad a = 1, 2.$$

Obviously, for the differentials, we have

$$\prod_{j=1}^W d\mathring{b}_{j,a} = \prod_{j=1}^W dc_{j,a}, \quad a = 1, 2.$$

In addition, for the domains, it is elementary to see

$$\mathring{\mathbf{b}}_a \in \mathring{\Upsilon} \iff \mathbf{c}_a \in \Upsilon, \quad a = 1, 2. \quad (8.39)$$

Now, we introduce the notation

$$\gamma_j^+ := \frac{1}{\sqrt{1 + a_+^2 + a_+^2 \lambda_j(S)}}, \quad \gamma_j^- := \frac{1}{\sqrt{1 + a_-^2 + a_-^2 \lambda_j(S)}},$$

and set the diagonal matrices

$$\mathbb{D}_+ := \text{diag}(\gamma_1^+, \dots, \gamma_W^+), \quad \mathbb{D}_- := \text{diag}(\gamma_1^-, \dots, \gamma_W^-).$$

By the assumption  $0 \leq E \leq \sqrt{2} - \kappa$  and (1.4), it is not difficult to check

$$|\gamma_j^+| \sim 1, \quad |\gamma_j^-| \sim 1, \quad \arg \gamma_j^+ \in \left(-\frac{\pi}{8}, 0\right], \quad \arg \gamma_j^- \in \left[0, \frac{\pi}{8}\right), \quad \forall j = 1, \dots, W. \quad (8.40)$$

With the notation introduced above, we have

$$\mathring{\mathbf{b}}'_1 \mathbb{A}_+ \mathring{\mathbf{b}}_1 = \mathbf{c}'_1 \mathbb{D}_+^{-2} \mathbf{c}_1, \quad \mathring{\mathbf{b}}'_2 \mathbb{A}_- \mathring{\mathbf{b}}_2 = \mathbf{c}'_2 \mathbb{D}_-^{-2} \mathbf{c}_2.$$

To simplify the following discussion, we enlarge the domain of the  $\mathbf{c}$ -variables to

$$\mathbf{c}_a \in \Upsilon_\infty \equiv \Upsilon_\infty(\varepsilon) := [-\Theta^{\frac{1}{2}}, \Theta^{\frac{1}{2}}]^W, \quad a = 1, 2.$$

Obviously,  $\mathring{\Upsilon} \subset \Upsilon_\infty$ . It is easy to check that (7.26) also holds when  $\mathbf{c}_a \in \Upsilon_\infty \setminus \mathring{\Upsilon}$  for either  $a = 1$  or  $2$ , according to (8.39), thus such a modification of the domain will only produce an error term of order  $O(\exp\{-\Theta\})$  in the integral (8.32), by using (8.7).

Now we do the scaling

$$\mathbf{c}_1 \rightarrow \mathbb{D}_+ \mathbf{c}_1, \quad \mathbf{c}_2 \rightarrow \mathbb{D}_- \mathbf{c}_2.$$

Consequently, we have

$$\mathring{\mathbf{b}}_1 = \mathbb{U}\mathbb{D}_+\mathbf{c}_1, \quad \mathring{\mathbf{b}}_2 = \mathbb{U}\mathbb{D}_-\mathbf{c}_2, \quad (8.41)$$

thus

$$\mathring{\mathbf{b}}_1'\mathbb{A}_+\mathring{\mathbf{b}}_1 = \sum_{j=1}^W c_{j,1}^2, \quad \mathring{\mathbf{b}}_2'\mathbb{A}_-\mathring{\mathbf{b}}_2 = \sum_{j=1}^W c_{j,2}^2. \quad (8.42)$$

Accordingly, we should adjust the change of differentials as

$$\prod_{j=1}^W d\mathring{b}_{j,1} \rightarrow \det \mathbb{D}_+ \cdot \prod_{j=1}^W dc_{j,1}, \quad \prod_{j=1}^W d\mathring{b}_{j,2} \rightarrow \det \mathbb{D}_- \cdot \prod_{j=1}^W dc_{j,2}.$$

In addition, the domain of  $\mathbf{c}_1$  should be changed from  $\Upsilon_\infty$  to  $\prod_{j=1}^W \mathbb{J}_j^+$ , where

$$\mathbb{J}_j^+ := (\gamma_j^+)^{-1}[-\Theta^{\frac{1}{2}}, \Theta^{\frac{1}{2}}],$$

and that of  $\mathbf{c}_2$  should be changed from  $\Upsilon_\infty$  to  $\prod_{j=1}^W \mathbb{J}_j^-$ , where

$$\mathbb{J}_j^- := (\gamma_j^-)^{-1}[-\Theta^{\frac{1}{2}}, \Theta^{\frac{1}{2}}].$$

By the fact  $\det \mathbb{D}_+\mathbb{D}_- = 1/\sqrt{\det \mathbb{A}_+\mathbb{A}_-}$ , we can write (8.32) as

$$\begin{aligned} 2^W \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \Upsilon_S) &= \frac{M^2}{(n!)^2 4^W \pi^{4W+2}} \cdot \frac{1}{\sqrt{\det \mathbb{A}_+\mathbb{A}_-}} \\ &\times \int_{\prod_{j=1}^W \mathbb{J}_j^+} \prod_{j=1}^W dc_{j,1} \int_{\prod_{j=1}^W \mathbb{J}_j^-} \prod_{j=1}^W dc_{j,2} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\hat{x}_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\hat{x}_{j,2} \\ &\times \int_{\hat{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,1} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,2} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W dv_{j,1} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W dv_{j,2} \prod_{j=1}^W \exp\left\{i \frac{\hat{x}_{j,1} + \hat{x}_{j,2}}{\sqrt{M}}\right\} \\ &\times \exp\left\{-\frac{1}{2} \sum c_{j,1}^2 - \frac{1}{2} \sum c_{j,2}^2 - R^b\right\} \cdot \exp\left\{-\frac{1}{2} \hat{\mathbf{x}}_1' \mathbb{A}_+ \hat{\mathbf{x}}_1 - \frac{1}{2} \hat{\mathbf{x}}_2' \mathbb{A}_- \hat{\mathbf{x}}_2 - R_\pm^x\right\} \\ &\times \exp\left\{(a_+ - a_-)^2 \sum_{a=1,2} \tau_a' S^{(1)} \tau_a - R^{t,b}\right\} \cdot \exp\left\{(a_+ - a_-)^2 \sum_{a=1,2} \mathbf{v}_a' S^{(1)} \mathbf{v}_a - R_\pm^{v,x}\right\} \\ &\times \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot A(\hat{X}, \hat{B}, V, T) + O(e^{-\Theta}). \end{aligned} \quad (8.43)$$

In (8.43), all the  $\mathbf{b}$  and  $\mathring{\mathbf{b}}$ -variables in the integrand should be regarded as the functions of the  $\mathbf{c}$ -variables, see (8.41).

Now, we consider the integrand as a function of  $\mathbf{c}$ -variables on the disks, namely,

$$c_{j,1} \in \mathbb{O}_j^+ := \{z \in \mathbb{C} : |z| \leq \Theta^{\frac{1}{2}} |\gamma_j^+|^{-1}\}, \quad c_{j,2} \in \mathbb{O}_j^- := \{z \in \mathbb{C} : |z| \leq \Theta^{\frac{1}{2}} |\gamma_j^-|^{-1}\}.$$

For  $\mathbf{c}_1 \in \prod_{j=1}^W \mathbb{O}_j^+$  and  $\mathbf{c}_2 \in \prod_{j=1}^W \mathbb{O}_j^-$ , by (8.40) and (8.41) we have

$$\|\mathring{\mathbf{b}}_1\|_\infty, \|\mathring{\mathbf{b}}_2\|_\infty \leq O(\Theta). \quad (8.44)$$

Here we used the elementary fact  $\|U\mathbf{a}\|_\infty \leq \sqrt{W}\|\mathbf{a}\|_\infty$  for any  $\mathbf{a} \in \mathbb{C}^W$  and unitary matrix  $U$ . Then, we deform the contour of  $c_{j,1}$  from  $\mathbb{J}_j^+$  to

$$(-\Sigma_j^+) \cup \mathbb{L}_j^+ \cup \Sigma_j^+$$

for each  $j = 1, \dots, W$ , where

$$\mathbb{L}_j^+ := \mathbb{R} \cap \mathbb{O}_j^+, \quad \Sigma_j^+ = \{z \in \partial \mathbb{O}_j^+ : 0 \leq \arg z \leq -\arg \gamma_j^+\}.$$



It is not difficult to see that

$$\text{Rec}_{j,1}^2 \geq \Theta, \quad \text{for } c_{j,1} \in (-\Sigma_j^+) \cup \Sigma_j^+,$$

by using (8.40). Consequently, by (8.42), we have

$$\left| \exp \left\{ -\frac{1}{2} \mathring{\mathbf{b}}_1' \mathbb{A}_+ \mathring{\mathbf{b}}_1 \right\} \right| = \left| \exp \left\{ -\frac{1}{2} \sum_{j=1}^W c_{j,1}^2 \right\} \right| \leq O(e^{-\Theta}).$$

Then using (8.7), we can get rid of the integral over  $\Sigma_j^+$  and  $-\Sigma_j^+$ , analogously to the discussion in Section 7. Similarly, we can perform the same argument for  $\mathbf{c}_2$ . Consequently, we can restrict the integral in (8.43) to the domain

$$\mathbf{c}_1 \in \prod_{j=1}^W \mathbb{L}_j^+, \quad \mathbf{c}_2 \in \prod_{j=1}^W \mathbb{L}_j^-.$$

So we can assume that  $\prod_{j=1}^W \mathbb{J}_j^+$  and  $\prod_{j=1}^W \mathbb{J}_j^-$  are replaced with  $\prod_{j=1}^W \mathbb{L}_j^+$  and  $\prod_{j=1}^W \mathbb{L}_j^-$  respectively in (8.43).

By (8.15), (8.44) and the fact  $\|\mathbf{a}\|_3^3 \leq \|\mathbf{a}\|_\infty \|\mathbf{a}\|_2^2$  for any vector  $\mathbf{a}$ , we see that

$$|R^b| \leq C \frac{\|\mathring{\mathbf{b}}_1\|_\infty + \|\mathring{\mathbf{b}}_2\|_\infty}{\sqrt{M}} (\|\mathring{\mathbf{b}}_1\|_2^2 + \|\mathring{\mathbf{b}}_2\|_2^2) \leq \frac{\Theta}{\sqrt{M}} (\|\mathbf{c}_1\|_2^2 + \|\mathbf{c}_2\|_2^2) \quad (8.45)$$

for some positive constant  $C$ , where in the last step we also used the fact that  $\|\mathbf{b}_a\|_2 = O(\|\mathbf{c}_a\|_2)$  for  $a = 1, 2$ , which is implied by (8.41) and (8.40). Consequently, we have

$$\exp \left\{ -\frac{1}{2} \|\mathbf{c}_1\|_2^2 - \frac{1}{2} \|\mathbf{c}_2\|_2^2 - R^b \right\} = \exp \left\{ -\left(\frac{1}{2} + o(1)\right) \|\mathbf{c}_1\|_2^2 - \left(\frac{1}{2} + o(1)\right) \|\mathbf{c}_2\|_2^2 \right\}.$$

This allows us to take a step further to truncate  $\mathbf{c}_1$  and  $\mathbf{c}_2$  according to their 2-norm, namely

$$\mathbf{c}_1, \mathbf{c}_2 \in \mathring{\Upsilon}. \quad (8.46)$$

Similarly to the discussion in the proof of Lemma 7.7, such a truncation will only produce an error of order  $\exp\{-\Theta\}$  to the integral, by using (8.7).

Now, analogously to (8.41), we can change  $\mathring{\mathbf{x}}$ -variables to  $\mathbf{d}$ -variables, defined by

$$\mathbf{d}_1 = (d_{1,1}, \dots, d_{W,1}) := \mathbb{D}_+^{-1} \mathbf{U}' \mathring{\mathbf{x}}_1, \quad \mathbf{d}_2 = (d_{1,2}, \dots, d_{W,2}) := \mathbb{D}_-^{-1} \mathbf{U}' \mathring{\mathbf{x}}_2.$$

Thus accordingly, we change the differentials

$$\prod_{j=1}^W d\mathring{x}_{j,1} \rightarrow \det \mathbb{D}_+ \cdot \prod_{j=1}^W dd_{j,1}, \quad \prod_{j=1}^W d\mathring{x}_{j,2} \rightarrow \det \mathbb{D}_- \cdot \prod_{j=1}^W dd_{j,2}.$$

In addition, like (8.46), we deform the domain to

$$\mathbf{d}_1, \mathbf{d}_2 \in \mathring{\Upsilon}.$$

Finally, from (8.43), we arrive at the representation

$$\begin{aligned} & 2^W \mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \Upsilon_S) \\ &= \frac{M^2}{(n!)^{24W} \pi^{4W+2}} \cdot \frac{1}{\det \mathbb{A}_+ \mathbb{A}_-} \cdot \int_{\mathring{\Upsilon}} \prod_{j=1}^W dc_{j,1} \int_{\mathring{\Upsilon}} \prod_{j=1}^W dc_{j,2} \int_{\mathring{\Upsilon}} \prod_{j=1}^W dd_{j,1} \int_{\mathring{\Upsilon}} \prod_{j=1}^W dd_{j,2} \\ & \times \int_{\mathring{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,1} \int_{\mathring{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,2} \int_{\mathring{\Upsilon}_S} \prod_{j=2}^W dv_{j,1} \int_{\mathring{\Upsilon}_S} \prod_{j=2}^W dv_{j,2} \prod_{j=1}^W \exp \left\{ i \frac{\mathring{x}_{j,1} + \mathring{x}_{j,2}}{\sqrt{M}} \right\} \\ & \times \exp \left\{ -\frac{1}{2} \|\mathbf{c}_1\|_2^2 - \frac{1}{2} \|\mathbf{c}_2\|_2^2 - R^b \right\} \cdot \exp \left\{ -\frac{1}{2} \|\mathbf{d}_1\|_2^2 - \frac{1}{2} \|\mathbf{d}_2\|_2^2 - R_\pm^x \right\} \end{aligned}$$

$$\begin{aligned}
& \times \exp \left\{ (a_+ - a_-)^2 \sum_{a=1,2} \tau'_a S^{(1)} \tau_a - R^{t,b} \right\} \cdot \exp \left\{ (a_+ - a_-)^2 \sum_{a=1,2} \mathbf{v}'_a S^{(1)} \mathbf{v}_a - R_{\pm}^{v,x} \right\} \\
& \times \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot \mathbf{A}(\hat{X}, \hat{B}, V, T) + O(e^{-\Theta}),
\end{aligned} \tag{8.47}$$

in which  $\mathbf{x}$  and  $\hat{\mathbf{x}}$ -variables should be regarded as functions of the  $\mathbf{d}$ -variables, as well,  $\mathbf{b}$  and  $\hat{\mathbf{b}}$ -variables should be regarded as functions of the  $\mathbf{c}$ -variables.

Now, in the Type II and III vicinities, we only do the change of coordinates for the  $\hat{\mathbf{b}}$ -variables, which is enough for our purpose. Consequently, we have

$$\begin{aligned}
\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_{\varkappa}^x, \Upsilon_{\varkappa}^x, \Upsilon_S, \mathbb{I}^{W-1}) &= \exp \{ M(K(D_{\pm}, I) - K(D_{\varkappa}, I)) \} \cdot \frac{(-a_{\varkappa}^2)^W}{(n!)^2} \cdot \frac{M^{W+1}}{8^W \pi^{3W+3}} \\
&\times \frac{1}{\sqrt{\det \mathbb{A}_+ \mathbb{A}_-}} \cdot \int_{\mathbb{I}^{W-1}} \prod_{j=2}^W \frac{d\theta_j}{2\pi} \int_{\mathbb{I}^{W-1}} \prod_{j=2}^W 2v_j dv_j \int_{\hat{\Upsilon}} \prod_{j=1}^W dc_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W dc_{j,2} \\
&\times \int_{\hat{\Upsilon}} \prod_{j=1}^W d\hat{x}_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W d\hat{x}_{j,2} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,1} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,2} \prod_{j=1}^W \exp \left\{ i \frac{\hat{x}_{j,1} + \hat{x}_{j,2}}{\sqrt{M}} \right\} \\
&\times \exp \left\{ -\frac{1}{2} \|\mathbf{d}_1\|_2^2 - \frac{1}{2} \|\mathbf{d}_2\|_2^2 - R^b \right\} \cdot \exp \left\{ (a_+ - a_-)^2 \sum_{a=1,2} \tau'_a S^{(1)} \tau_a - R^{t,b} \right\} \\
&\times \exp \left\{ -\frac{1}{2} \hat{\mathbf{x}}' \mathbb{A}_{\varkappa}^v \hat{\mathbf{x}} - R_{\varkappa}^x - R_{\varkappa}^{v,x} \right\} \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 \cdot \mathbf{A}(\hat{X}, \hat{B}, V, T) + O(e^{-\Theta}).
\end{aligned} \tag{8.48}$$

By (8.40) and (8.41), it is easy to see

$$\mathbf{c}_1, \mathbf{c}_2 \in \hat{\Upsilon} \implies \hat{\mathbf{b}}_1, \hat{\mathbf{b}}_2 \in \hat{\Upsilon} \tag{8.49}$$

Similarly, we have

$$\mathbf{d}_1, \mathbf{d}_2 \in \hat{\Upsilon} \implies \hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2 \in \hat{\Upsilon}. \tag{8.50}$$

We keep the terminology “Type I’, II and III vicinities” for the slightly modified domains defined in terms of  $\mathbf{c}$ ,  $\mathbf{d}$ ,  $\boldsymbol{\tau}$  and  $\mathbf{v}$ -variables. More specifically, we redefine the vicinities as follows.

**Definition 8.5.** *We slightly modify Definition 8.1 as follows.*

- *Type I’ vicinity:*  $\mathbf{c}_1, \mathbf{c}_2, \mathbf{d}_1, \mathbf{d}_2 \in \hat{\Upsilon}, \quad \boldsymbol{\tau}_1, \boldsymbol{\tau}_2, \mathbf{v}_1, \mathbf{v}_2 \in \hat{\Upsilon}_S.$
- *Type II vicinity:*  $\mathbf{c}_1, \mathbf{c}_2, \hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2 \in \hat{\Upsilon}, \quad \boldsymbol{\tau}_1, \boldsymbol{\tau}_2 \in \hat{\Upsilon}_S, \quad V_j \in \hat{U}(2) \text{ for all } j = 2, \dots, W,$   
where  $\hat{\mathbf{x}}$ -variables are defined in (8.8) with  $\varkappa = +$ .
- *Type III vicinity:*  $\mathbf{c}_1, \mathbf{c}_2, \hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2 \in \hat{\Upsilon}, \quad \boldsymbol{\tau}_1, \boldsymbol{\tau}_2 \in \hat{\Upsilon}_S, \quad V_j \in \hat{U}(2) \text{ for all } j = 2, \dots, W,$   
where  $\hat{\mathbf{x}}$ -variables are defined in (8.8) with  $\varkappa = -$ .

Now, recall the remainder terms  $R^b$ ,  $R_{\pm}^x$ ,  $R_+^x$  and  $R_-^x$  in Lemma 8.2,  $R^{t,b}$  and  $R_{\pm}^{v,x}$  in (8.23),  $R_+^{v,x}$  in (8.33) and  $R_-^{v,x}$  in (8.36). In light of (8.49) and (8.50), the bounds on these remainder terms are the same as those obtained in Section 8.1. For the convenience of the reader, we collect them as the following proposition.

**Proposition 8.6.** *Under Assumptions 1.1 and 1.13, we have the following estimate, in the vicinities.*

$$(i): \quad R^{t,b} = O\left(\frac{\Theta^4}{M}\right), \quad R_{\pm}^{v,x} = O\left(\frac{\Theta^4}{M}\right),$$

$$\begin{aligned}
(ii): \quad R^b &= O\left(\frac{\Theta^2}{\sqrt{M}}\right), \quad R^x = O\left(\frac{\Theta^2}{\sqrt{M}}\right), \quad R_+^x = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right), \quad R_-^x = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right), \\
(iii): \quad R_+^{v,x} &= O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right), \quad R_-^{v,x} = O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right).
\end{aligned}$$

*Proof.* Note that, (i) can be obtained from (8.26), and (ii) follows from Lemma 8.2, and (iii) is implied by (8.34) and (8.37). Hence, we completed the proof.  $\square$

Analogously, in the vicinities,  $\|\mathring{\mathbf{b}}_1\|_2^2$ ,  $\|\mathring{\mathbf{b}}_2\|_2^2$ ,  $\|\mathring{\mathbf{x}}_1\|_2^2$ ,  $\|\mathring{\mathbf{x}}_2\|_2^2$ ,  $\|\mathring{\mathbf{t}}\|_\infty$  and  $\|\mathring{\mathbf{v}}\|_\infty$  are still bounded by  $\Theta$ .

## 9. INTEGRAL OVER THE TYPE I VICINITIES

With (8.47), we estimate the integral over the Type I vicinity in this section. At first, in the Type I' vicinity, we have  $\|\mathring{\mathbf{x}}_a\|_\infty = O(\Theta^{\frac{1}{2}})$  and  $\|\mathring{\mathbf{b}}_a\|_\infty = O(\Theta^{\frac{1}{2}})$  for  $a = 1, 2$ . Consequently, according to the parametrization in (8.2), we have

$$x_{j,1} - x_{j,2} = a_+ - a_- + O\left(\frac{\Theta^{\frac{1}{2}}}{\sqrt{M}}\right), \quad b_{j,1} + b_{j,2} = a_+ - a_- + O\left(\frac{\Theta^{\frac{1}{2}}}{\sqrt{M}}\right), \quad (9.1)$$

which implies

$$\prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 = (a_+ - a_-)^{4W} \left(1 + O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right)\right). \quad (9.2)$$

Hence, what remains is to estimate the function  $A(\hat{X}, \hat{B}, V, T)$ . We have the following lemma.

**Lemma 9.1.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type I' vicinity, for any given positive integer  $n$ , there is  $N_0 = N_0(n)$ , such that for all  $N \geq N_0$  we have*

$$|A(\hat{X}, \hat{B}, V, T)| \leq \frac{\Theta^2 W^{C_0}}{M(N\eta)^{n+\ell}} \cdot |\det \mathbb{A}_+|^2 \cdot \det(S^{(1)})^2.$$

for some positive constant  $C_0$  and some integer  $\ell = O(1)$ , both of which are independent of  $n$ .

With (8.47), (9.2) and Lemma 9.1, we can prove Lemma 5.8.

*Proof of Lemma 5.8.* Using (8.47), (9.2), Lemma 9.1, Proposition 8.6 with (5.38), the fact  $\det \mathbb{A}_+ = \det \mathbb{A}_-$  and the trivial estimate

$$M\Theta^2 W^{C_0} \frac{1}{(N\eta)^\ell} \leq N^{C_0}$$

for sufficiently large constant  $C_0$ , we have

$$\begin{aligned}
2^W |\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \Upsilon_S)| &\leq \frac{N^{C_0}}{(N\eta)^n} \cdot \frac{1}{(2\pi^2)^{2W}} \cdot \det(S^{(1)})^2 \cdot (a_+ - a_-)^{4W} \\
&\times \int_{\hat{\Upsilon}} \prod_{j=1}^W dc_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W dc_{j,2} \int_{\hat{\Upsilon}} \prod_{j=1}^W dd_{j,1} \int_{\hat{\Upsilon}} \prod_{j=1}^W dd_{j,2} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,1} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W d\tau_{j,2} \\
&\times \int_{\hat{\Upsilon}_S} \prod_{j=2}^W dv_{j,1} \int_{\hat{\Upsilon}_S} \prod_{j=2}^W dv_{j,2} \exp \left\{ -\frac{1}{2} (\|\mathbf{c}_1\|_2^2 + \|\mathbf{c}_2\|_2^2 + \|\mathbf{d}_1\|_2^2 + \|\mathbf{d}_2\|_2^2) \right\} \\
&\times \exp \left\{ (a_+ - a_-)^2 \left( \boldsymbol{\tau}'_1 S^{(1)} \boldsymbol{\tau}_1 + \boldsymbol{\tau}'_1 S^{(1)} \boldsymbol{\tau}_2 + \mathbf{v}'_1 S^{(1)} \mathbf{v}_1 + \mathbf{v}'_1 S^{(1)} \mathbf{v}_2 \right) \right\} + O(e^{-\Theta}).
\end{aligned}$$

Then, by elementary Gaussian integral we obtain (5.40). Hence, we completed the proof of Lemma 5.8.  $\square$

The remaining part of this section will be dedicated to the proof of Lemma 9.1. Recall the definitions of the functions  $A(\cdot)$ ,  $Q(\cdot)$ ,  $P(\cdot)$  and  $F(\cdot)$  in (3.31), (4.3), (4.4) and (4.5). Using the strategy in Section 6 again, we ignore the irrelevant factor  $Q(\cdot)$  at the beginning. Hence, we bound  $P(\cdot)$  and  $F(\cdot)$  at first, and modify the bounding procedure slightly to take  $Q(\cdot)$  into account in the end, resulting a proof of Lemma 9.1.

**9.1.  $P(\hat{X}, \hat{B}, V, T)$  in the Type I' vicinity.** As mentioned above, we should always regard  $\mathbf{b}$  or  $\mathring{\mathbf{b}}$ -variables as functions of  $\mathbf{c}$ -variables, regard  $\mathbf{x}$  or  $\mathring{\mathbf{x}}$ -variables as functions of  $\mathbf{d}$ -variables. Our aim, in this section, is to prove the following lemma.

**Lemma 9.2.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type I' vicinity, we have*

$$P(\hat{X}, \hat{B}, V, T) \leq \frac{W^{2+\gamma}\Theta^2}{M} |\det \mathbb{A}_+|^2 \det(S^{(1)})^2. \quad (9.3)$$

Before commencing the formal proof, we introduce more notation below. In the sequel, we will use the notation

$$\mathring{\kappa}_j \equiv \mathring{\kappa}_j(\hat{X}, \hat{B}, V, T) := |\mathring{x}_{j,1}| + |\mathring{x}_{j,2}| + |\mathring{b}_{j,1}| + |\mathring{b}_{j,2}| + |\mathring{v}_j| + |\mathring{t}_j| = O(\Theta), \quad (9.4)$$

where the bound holds in the Type I' vicinity, according to the facts  $\|\mathring{\mathbf{x}}_a\|_\infty = O(\Theta^{\frac{1}{2}})$ ,  $\|\mathring{\mathbf{b}}_a\|_\infty = O(\Theta^{\frac{1}{2}})$  for  $a = 1, 2$ ,  $\|\mathring{\mathbf{t}}\|_\infty = O(\Theta)$  and  $\|\mathring{\mathbf{v}}\|_\infty = O(\Theta)$ .

Recalling (6.17) with  $\varpi_j$  defined in (6.13) and  $\hat{\varpi}_j$  in (6.15). Now, we write

$$\begin{aligned} \varpi_j &= \exp \left\{ -M \log \det (1 + M^{-1} V_j^* \hat{X}_j^{-1} V_j \Omega_j T_j^{-1} \hat{B}_j^{-1} T_j \Xi_j) \right\} \\ &=: \exp \left\{ -\text{Tr} V_j^* \hat{X}_j^{-1} V_j \Omega_j T_j^{-1} \hat{B}_j^{-1} T_j \Xi_j \right\} \exp \left\{ \sum_{\ell=2}^4 \frac{(-1)^{\ell-1}}{\ell M^{\ell-1}} \Delta_{\ell,j} \right\}, \end{aligned} \quad (9.5)$$

where

$$\Delta_{\ell,j} := \text{Tr} (V_j^* \hat{X}_j^{-1} V_j \Omega_j T_j^{-1} \hat{B}_j^{-1} T_j \Xi_j)^\ell. \quad (9.6)$$

The second step of (9.5) follows from the Taylor expansion of the logarithmic function. The expansion terminates at the 4th order term since  $\Delta_{\ell,j}$  is a homogeneous polynomial of  $\Omega_j$  and  $\Xi_j$ -variables with degree  $2\ell$ , regarding all the complex variables as fixed parameters. Now, we expand the first factor of (9.5) around the Type I' saddle point, namely

$$\exp \left\{ -\text{Tr} V_j^* \hat{X}_j^{-1} V_j \Omega_j T_j^{-1} \hat{B}_j^{-1} T_j \Xi_j \right\} =: \exp \left\{ -\text{Tr} D_\pm^{-1} \Omega_j D_\pm^{-1} \Xi_j \right\} \exp \left\{ -\frac{1}{\sqrt{M}} \Delta_j \right\}. \quad (9.7)$$

We take (9.7) as the definition of  $\Delta_j$ , which is of the form

$$\Delta_j = \sum_{\alpha, \beta=1}^4 \mathring{\mathbf{p}}_{j, \alpha, \beta} \cdot \omega_{j, \alpha} \xi_{j, \beta}$$

for some function  $\mathring{\mathbf{p}}_{j, \alpha, \beta}$  of  $\mathring{\mathbf{x}}$ ,  $\mathring{\mathbf{b}}$ ,  $\mathring{\mathbf{v}}$  and  $\mathring{\mathbf{t}}$ -variables, satisfying

$$\mathring{\mathbf{p}}_{j, \alpha, \beta} = O(\mathring{\kappa}_j), \quad \forall \alpha, \beta = 1, \dots, 4. \quad (9.8)$$

One can check (9.8) easily by using (8.8)-(8.10). Analogously, we can also write

$$\Delta_{\ell,j} = \sum_{\substack{\alpha_1, \dots, \alpha_\ell, \\ \beta_1, \dots, \beta_\ell=1}}^4 \mathring{\mathbf{p}}_{\ell,j, \alpha, \beta} \prod_{i=1}^{\ell} \omega_{j, \alpha_i} \xi_{j, \beta_i}, \quad \alpha := (\alpha_1, \dots, \alpha_\ell), \quad \beta := (\beta_1, \dots, \beta_\ell), \quad (9.9)$$

where

$$\mathring{\mathbf{p}}_{\ell,j, \alpha, \beta} = O(1), \quad \forall \ell = 2, \dots, 4; \alpha_1, \dots, \alpha_\ell, \beta_1, \dots, \beta_\ell = 1, \dots, 4. \quad (9.10)$$

The bound on  $\hat{\mathbf{p}}_{\ell,j,\alpha,\beta}$  in (9.10) follows from the fact that all the  $V_j$ ,  $\hat{X}_j^{-1}$ ,  $T_j$ ,  $T_j^{-1}$  and  $\hat{B}_j^{-1}$ -entries are bounded in the Type I' vicinity, uniformly in  $j$ . Consequently, we can write for  $j \neq p, q$

$$\exp \left\{ -\frac{1}{\sqrt{M}}\Delta_j + \sum_{\ell=2}^4 \frac{(-1)^{\ell-1}}{\ell M^{\ell-1}}\Delta_{\ell,j} \right\} = 1 + \sum_{\ell=1}^4 M^{-\frac{\ell}{2}} \sum_{\substack{\alpha_1, \dots, \alpha_\ell, \\ \beta_1, \dots, \beta_\ell=1}} \hat{\mathbf{q}}_{\ell,j,\alpha,\beta} \prod_{i=1}^{\ell} \omega_{j,\alpha_i} \xi_{j,\beta_i}. \quad (9.11)$$

In a similar manner, we can also write for  $k = p, q$ ,

$$\exp \left\{ -\frac{1}{\sqrt{M}}\Delta_k + \sum_{\ell=2}^4 \frac{(-1)^{\ell-1}}{\ell M^{\ell-1}}\Delta_{\ell,k} \right\} \hat{\omega}_k = \hat{\mathbf{p}}_0(\cdot) \left( 1 + \sum_{\ell=1}^4 M^{-\frac{\ell}{2}} \sum_{\substack{\alpha_1, \dots, \alpha_\ell, \\ \beta_1, \dots, \beta_\ell=1}} \hat{\mathbf{q}}_{\ell,k,\alpha,\beta} \prod_{i=1}^{\ell} \omega_{k,\alpha_i} \xi_{k,\beta_i} \right), \quad (9.12)$$

where  $\hat{\mathbf{p}}_0(\cdot) = \det \hat{X}_k / \det \hat{B}_k$ , which is introduced in (6.15), and  $\hat{\mathbf{q}}_{\ell,j,\alpha,\beta}$  is some function of  $\hat{X}$ ,  $\hat{B}$ ,  $V$  and  $T$ -variables, satisfying the bound

$$\hat{\mathbf{q}}_{\ell,j,\alpha,\beta} = O((1 + \hat{\kappa}_j)^\ell), \quad \forall \ell = 1, \dots, 4, \quad j = 1, \dots, W. \quad (9.13)$$

Obviously, we have  $\hat{\mathbf{p}}_0(\cdot) = O(1)$  in Type I' vicinity.

Now, in order to distinguish  $\ell, \alpha$  and  $\beta$  for different  $j$ , we index them as  $\ell_j, \alpha_j$  and  $\beta_j$ , where

$$\alpha_j \equiv \alpha_j(\ell_j) := (\alpha_{j,1}, \dots, \alpha_{j,\ell_j}), \quad \beta_j \equiv \beta_j(\ell_j) := (\beta_{j,1}, \dots, \beta_{j,\ell_j}).$$

In addition, we introduce the vector

$$\vec{\ell} := (\ell_1, \dots, \ell_W), \quad \vec{\alpha} \equiv \vec{\alpha}(\vec{\ell}) := (\alpha_1, \dots, \alpha_W), \quad \vec{\beta} \equiv \vec{\beta}(\vec{\ell}) := (\beta_1, \dots, \beta_W).$$

Let  $\|\vec{\ell}\|_1 = \sum_{j=1}^W \ell_j$  be the 1-norm of  $\vec{\ell}$ . Note that  $\vec{\alpha}$  and  $\vec{\beta}$  are  $\|\vec{\ell}\|_1$ -dimensional. With these notations, using (6.17), (9.5), (9.7), (9.11) and (9.12) we have the representation

$$\begin{aligned} \mathcal{P}(\Omega, \Xi, \hat{X}, \hat{B}, V, T) &= \hat{\mathbf{p}}_0(\hat{X}_p, \hat{B}_p) \hat{\mathbf{p}}_0(\hat{X}_q, \hat{B}_q) \cdot \exp \left\{ -\sum_{j,k} \hat{\mathbf{s}}_{jk} \text{Tr} \Omega_j \Xi_k - \sum_{j=1}^W \text{Tr} D_{\pm}^{-1} \Omega_j D_{\pm}^{-1} \Xi_j \right\} \\ &\times \left( 1 + \sum_{\substack{\vec{\ell} \in \llbracket 0, 4 \rrbracket^W, \\ \text{s.t. } \|\vec{\ell}\|_1 \geq 1}} M^{-\frac{\|\vec{\ell}\|_1}{2}} \sum_{\vec{\alpha}, \vec{\beta} \in \llbracket 1, 4 \rrbracket^{\|\vec{\ell}\|_1}} \prod_{j=1}^W \hat{\mathbf{q}}_{\ell_j, j, \alpha_j, \beta_j} \cdot \prod_{j=1}^W \prod_{i=1}^{\ell_j} \omega_{j, \alpha_{j,i}} \xi_{j, \beta_{j,i}} \right), \end{aligned} \quad (9.14)$$

where we made the convention

$$\hat{\mathbf{q}}_{0,j,\emptyset,\emptyset} = 1, \quad \prod_{i=1}^0 \omega_{j, \alpha_{j,i}} \xi_{j, \beta_{j,i}} = 1, \quad \forall j = 1, \dots, W. \quad (9.15)$$

According to (9.13) and (9.15), we have

$$\prod_{j=1}^W |\hat{\mathbf{q}}_{\ell_j, j, \alpha_j, \beta_j}| \leq e^{O(\|\vec{\ell}\|_1)} \prod_{j=1}^W (1 + \hat{\kappa}_j)^{\ell_j}. \quad (9.16)$$

In addition, we can decompose the sum

$$\sum_{\substack{\vec{\ell} \in \llbracket 0, 4 \rrbracket^W, \\ \text{s.t. } \|\vec{\ell}\|_1 \geq 1}} = \sum_{\mathbf{m}=1}^{4W} \sum_{\substack{\vec{\ell} \in \llbracket 0, 4 \rrbracket^W, \\ \text{s.t. } \|\vec{\ell}\|_1 = \mathbf{m}}} . \quad (9.17)$$

It is easy to see

$$\#\{\vec{\ell} \in \llbracket 0, 4 \rrbracket^W : \|\vec{\ell}\|_1 = \mathbf{m}\} \leq \binom{4W}{\mathbf{m}}. \quad (9.18)$$

Moreover, it is obvious that

$$\#\{\vec{\alpha}, \vec{\beta} \in \llbracket 1, 4 \rrbracket^{\|\vec{\ell}\|_1}\} = 16^{\|\vec{\ell}\|_1}. \quad (9.19)$$

Therefore, it suffices to investigate the integral

$$\mathfrak{P}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}} := \int d\Omega d\Xi \exp \left\{ - \sum_{j,k} \tilde{s}_{jk} \text{Tr} \Omega_j \Xi_k - \sum_{j=1}^W \text{Tr} D_{\pm}^{-1} \Omega_j D_{\pm}^{-1} \Xi_j \right\} \prod_{j=1}^W \prod_{i=1}^{\ell_j} \omega_{j, \alpha_{j,i}} \xi_{j, \beta_{j,i}}$$

for each combination  $(\vec{\ell}, \vec{\alpha}, \vec{\beta})$ , and then sum it up for  $(\vec{\ell}, \vec{\alpha}, \vec{\beta})$  to get the estimate of  $P(\hat{X}, \hat{B}, V, T)$ . Specifically, we have the following lemma.

**Lemma 9.3.** *With the notation above, we have*

$$\mathfrak{P}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}} = 0, \quad \text{if } \|\vec{\ell}\|_1 = 0 \quad \text{or} \quad 1. \quad (9.20)$$

Moreover, we have

$$|\mathfrak{P}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}}| \leq |\det \mathbb{A}_+|^2 \det(S^{(1)})^2 (\|\vec{\ell}\|_1 - 1)! (2W^\gamma)^{(\|\vec{\ell}\|_1 - 1)}, \quad \text{if } \|\vec{\ell}\|_1 \geq 2. \quad (9.21)$$

We postpone the proof of Lemma 9.3 and prove Lemma 9.2 at first.

*Proof of Lemma 9.2.* By (4.4), (9.14) and (9.20) and the fact that  $\hat{\mathbf{p}}_0(\cdot) = O(1)$ , we have

$$|P(\hat{X}, \hat{B}, V, T)| \leq C \sum_{\substack{\vec{\ell} \in \llbracket 0, 4 \rrbracket^W, \\ \text{s. t. } \|\vec{\ell}\|_1 \geq 2}} M^{-\frac{\|\vec{\ell}\|_1}{2}} \sum_{\vec{\alpha}, \vec{\beta} \in \llbracket 1, 4 \rrbracket^{\|\vec{\ell}\|_1}} \prod_{j=1}^W |\hat{\mathbf{q}}_{\ell_j, j, \alpha_j, \beta_j}| \cdot |\mathfrak{P}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}}|. \quad (9.22)$$

Substituting the bounds (9.16), (9.19) and (9.21) into (9.22) yields

$$\begin{aligned} |P(\hat{X}, \hat{B}, V, T)| &\leq |\det \mathbb{A}_+|^2 \det(S^{(1)})^2 \\ &\times \sum_{\substack{\vec{\ell} \in \llbracket 0, 4 \rrbracket^W, \\ \text{s. t. } \|\vec{\ell}\|_1 \geq 2}} e^{O(\|\vec{\ell}\|_1)} \cdot M^{-\frac{\|\vec{\ell}\|_1}{2}} \cdot (\|\vec{\ell}\|_1 - 1)! (2W^\gamma)^{(\|\vec{\ell}\|_1 - 1)} \cdot \prod_{j=1}^W (1 + \kappa_j)^{\ell_j}. \end{aligned} \quad (9.23)$$

Now, from (9.4) we have

$$\prod_{j=1}^W (1 + \kappa_j)^{\ell_j} \leq \Theta^{\|\vec{\ell}\|_1}, \quad (9.24)$$

which can absorb the irrelevant factor  $e^{O(\|\vec{\ell}\|_1)}$ . Using (9.17), (9.18) and (9.24), we have

$$\begin{aligned} &\sum_{\substack{\vec{\ell} \in \llbracket 0, 4 \rrbracket^W, \\ \text{s. t. } \|\vec{\ell}\|_1 \geq 2}} e^{O(\|\vec{\ell}\|_1)} \cdot M^{-\frac{\|\vec{\ell}\|_1}{2}} \cdot (\|\vec{\ell}\|_1 - 1)! (2W^\gamma)^{(\|\vec{\ell}\|_1 - 1)} \cdot \prod_{j=1}^W (1 + \kappa_j)^{\ell_j} \\ &\leq \sum_{\mathbf{m}=2}^{4W} \binom{4W}{\mathbf{m}} \cdot M^{-\frac{\mathbf{m}}{2}} \cdot \Theta^{\mathbf{m}} \cdot \mathbf{m}! W^{(\mathbf{m}-1)\gamma} \\ &\leq \sum_{\mathbf{m}=2}^{4W} (4W)^{\mathbf{m}} \cdot M^{-\frac{\mathbf{m}}{2}} \cdot \Theta^{\mathbf{m}} \cdot W^{(\mathbf{m}-1)\gamma} = O\left(\frac{W^{2+\gamma} \Theta^2}{M}\right), \end{aligned} \quad (9.25)$$

where the last step follows from (5.33) and (5.38). Now, substituting (9.25) into (9.23), we can complete the proof of Lemma 9.2.  $\square$

Hence, what remains is to prove Lemma 9.3. We will need the following technical lemma whose proof is postponed.

**Lemma 9.4.** *For any index sets  $\mathbf{I}, \mathbf{J} \subset \{1, \dots, W\}$  with  $|\mathbf{I}| = |\mathbf{J}| = \mathbf{m} \geq 1$ , we have the following bounds for the determinants of the submatrices of  $S$ ,  $\mathbb{A}_+$  and  $\mathbb{A}_-$  defined in (8.14).*

- For  $(\mathbb{A}_+)^{(\mathbf{l}|\mathbf{J})}$  and  $(\mathbb{A}_-)^{(\mathbf{l}|\mathbf{J})}$ , we have

$$\frac{|\det(\mathbb{A}_+)^{(\mathbf{l}|\mathbf{J})}|}{|\det \mathbb{A}_+|} \leq 1, \quad \frac{|\det(\mathbb{A}_-)^{(\mathbf{l}|\mathbf{J})}|}{|\det \mathbb{A}_-|} \leq 1. \quad (9.26)$$

- For  $S^{(\mathbf{l}|\mathbf{J})}$ , we have

$$\frac{|\det S^{(\mathbf{l}|\mathbf{J})}|}{|\det S^{(1)}|} \leq (\mathfrak{m}-1)!(2W^\gamma)^{(\mathfrak{m}-1)}. \quad (9.27)$$

*Proof of Lemma 9.3.* Recall the definition in (6.20). Furthermore, we introduce the matrix

$$\mathbb{H} = (a_+^{-2} \mathbb{A}_+) \oplus S \oplus S \oplus (a_-^{-2} \mathbb{A}_-). \quad (9.28)$$

Using the fact  $a_+ a_- = -1$ , we can write

$$-\sum_{j,k} \tilde{\mathfrak{s}}_{jk} \text{Tr} \Omega_j \Xi_k - \sum_{j=1}^W \text{Tr} D_\pm^{-1} \Omega_j D_\pm^{-1} \Xi_j = -\vec{\Omega} \mathbb{H} \vec{\Xi}'.$$

Now, by the Gaussian integral of the Grassmann variables in (3.2), we see that

$$|\mathfrak{P}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}}| = |\det \mathbb{H}^{(\mathbf{l}|\mathbf{J})}| \quad (9.29)$$

for some index sets  $\mathbf{l}, \mathbf{J} \subset \{1, \dots, 4W\}$  determined by  $\vec{\alpha}$  and  $\vec{\beta}$  such that

$$|\mathbf{l}| = |\mathbf{J}| = \|\vec{\ell}\|_1.$$

Here we mention that (9.29) may fail when at least two components in  $\alpha_j$  coincide for some  $j$ . But  $\mathfrak{P}_{\vec{\ell}, \vec{\alpha}, \vec{\beta}} = 0$  in this case because of  $\chi^2 = 0$  for any Grassmann variable  $\chi$ .

Now, obviously, there exists index sets  $\mathbf{l}_\alpha, \mathbf{J}_\alpha \subset \{1, \dots, W\}$  for  $\alpha = 1, \dots, 4$  such that

$$\mathbb{H}^{(\mathbf{l}|\mathbf{J})} = (a_+^{-2} \mathbb{A}_+)^{(\mathbf{l}_1|\mathbf{J}_1)} \oplus S^{(\mathbf{l}_2|\mathbf{J}_2)} \oplus S^{(\mathbf{l}_3|\mathbf{J}_3)} \oplus (a_-^{-2} \mathbb{A}_-)^{(\mathbf{l}_4|\mathbf{J}_4)}, \quad \sum_\alpha |\mathbf{l}_\alpha| = \sum_\alpha |\mathbf{J}_\alpha| = \|\vec{\ell}\|_1.$$

It suffices to consider the case

$$|\mathbf{l}_\alpha| = |\mathbf{J}_\alpha|, \quad \forall \alpha = 1, 2, 3, 4.$$

Otherwise,  $\det \mathbb{H}^{(\mathbf{l}|\mathbf{J})}$  is obviously 0, in light of the block structure of  $\mathbb{H}$ , see the definition (9.28). Now, note that, since  $\det S = 0$ , we have

$$\det \mathbb{H}^{(\mathbf{l}|\mathbf{J})} = 0, \quad \text{if } \|\vec{\ell}\|_1 = 0, 1.$$

In addition, if  $\|\vec{\ell}\|_1 = 2$ , by using (9.30) below, one has

$$|\det \mathbb{H}^{(\mathbf{l}|\mathbf{J})}| = \begin{cases} |\det \mathbb{A}_+ \mathbb{A}_-| \det(S^{(1)})^2, & \text{if } |\mathbf{l}_2| = |\mathbf{l}_3| = |\mathbf{J}_2| = |\mathbf{J}_3| = 1, \\ 0, & \text{otherwise.} \end{cases}$$

For more general  $\vec{\ell}$ , by Lemma 9.4, we have

$$|\det \mathbb{H}^{(\mathbf{l}|\mathbf{J})}| \leq |\det \mathbb{A}_+ \mathbb{A}_-| \det(S^{(1)})^2 (\|\vec{\ell}\|_1 - 1)!(2W^\gamma)^{(\|\vec{\ell}\|_1 - 1)}.$$

Then, by the fact  $|\det \mathbb{A}_+ \mathbb{A}_-| = |\det \mathbb{A}_+|^2$ , we can conclude the proof of Lemma 9.3.  $\square$

To prove Lemma 9.4, we will need the following lemma.

**Lemma 9.5.** *For the weighted Laplacian  $S$ , we have*

$$\det S^{(i|j)} = (-1)^{j-i} \det S^{(i)}, \quad \forall i, j = 1, \dots, W \quad (9.30)$$

**Remark 9.6.** *A direct consequence of (9.30) is*

$$\det S^{(1)} = \dots = \det S^{(W)}. \quad (9.31)$$

*Proof of Lemma 9.5.* Without loss of generality, we assume  $j > i$  in the sequel. We introduce the matrices

$$P_{ij} := I_{i-1} \oplus \begin{pmatrix} I_{j-i-1} \\ 1 \end{pmatrix} \oplus I_{W-j}, \quad E_j := I - 2\mathbf{e}_j\mathbf{e}_j^* - \sum_{\ell \neq j} \mathbf{e}_\ell\mathbf{e}_\ell^*.$$

It is not difficult to check

$$S^{(i|j)} = S^{(i)} P_{ij} E_j. \quad (9.32)$$

Then, by the fact  $\det P_{ij} E_j = (-1)^{j-i}$ , we can get the conclusion.  $\square$

*Proof of Lemma 9.4.* At first, by the definition in (8.14), (1.4) and the fact  $\text{Re}a_+^2 = \text{Re}a_-^2 > 0$ , it is easy to see that the singular values of  $\mathbb{A}_+$  and  $\mathbb{A}_-$  are all larger than 1. With the aid of the rectangular matrix  $(\mathbb{A}_+)^{(|I|)}$  as an intermediate matrix, we can use Cauchy interlacing property twice to see that the  $k$ -th largest singular value of  $(\mathbb{A}_+)^{(|I|J)}$  is always smaller than the  $k$ -th largest singular value of  $\mathbb{A}_+$ . Consequently, we have the first inequality of (9.26). In the same manner, we can get the second inequality of (9.26).

Now, we prove (9.27). At first, we address the case that  $I \cap J \neq \emptyset$ . In light of (9.31), without loss of generality, we assume that  $1 \in I \cap J$ . Then  $S^{(|I|J)}$  is a submatrix of  $S^{(1)}$ . Therefore, we can find two permutation matrices  $P$  and  $Q$ , such that

$$PS^{(1)}Q = \begin{pmatrix} A & B \\ C & D \end{pmatrix},$$

where  $D = S^{(|I|J)}$ . Now, by Schur complement, we know that

$$\frac{|\det S^{(|I|J)}|}{|\det S^{(1)}|} = |\det(A - BD^{-1}C)^{-1}|.$$

Moreover,  $(A - BD^{-1}C)^{-1}$  is the  $(|I| - 1)$  by  $(|I| - 1)$  upper-left corner of

$$(PS^{(1)}Q)^{-1} = Q^{-1}(S^{(1)})^{-1}P^{-1}.$$

That means  $\det S^{(|I|J)} / \det S^{(1)}$  is the determinant of a sub matrix of  $(S^{(1)})^{-1}$  (with dimension  $|I| - 1$ ), up to a sign. Then, by Assumption 1.1 (iii), we can easily get

$$|\det S^{(|I|J)}| / |\det S^{(1)}| \leq (|I| - 1)! W^{(|I|-1)\gamma}.$$

Now, for the case  $I \cap J = \emptyset$ , we can fix one  $i \in I$  and  $j \in J$ . Due to (9.30), it suffices to consider

$$\frac{\det S^{(|I|J)}}{\det S^{(i|j)}}. \quad (9.33)$$

By similar discussion, one can see that (9.33) is the determinant of a sub matrix of  $(S^{(i|j)})^{-1}$  with dimension  $|I| - 1$ . Hence, it suffices to investigate the bound of the entries of  $(S^{(i|j)})^{-1}$ . From (9.32) we have

$$(S^{(i|j)})^{-1} = E_j^{-1} P_{ij}^{-1} (S^{(i)})^{-1}. \quad (9.34)$$

Observe that

$$P_{ij}^{-1} = I_{i-1} \oplus \begin{pmatrix} & 1 \\ I_{j-i-1} & \end{pmatrix} \oplus I_{W-j}, \quad E_j^{-1} = E_j.$$

Then, it is elementary to see that the entries of  $(S^{(i|j)})^{-1}$  are bounded by  $2W^\gamma$ , in light of (9.34) and Assumption 1.1 (iii). Consequently, we have

$$\frac{|\det S^{(|I|J)}|}{|\det S^{(i|j)}|} \leq (|I| - 1)! (2W^\gamma)^{|I|-1},$$

which implies (9.27). Hence, we completed the proof of Lemma 9.4.  $\square$



9.2.  $F(\hat{X}, \hat{B}, V, T)$  in the Type I' vicinity. Neglecting the  $X^{[1]}$ ,  $\mathbf{y}^{[1]}$  and  $\mathbf{w}^{[1]}$ -variables in  $Q(\cdot)$  at first, we investigate the integral  $F(\hat{X}, \hat{B}, V, T)$  in the Type I' vicinity in this section. We have the following lemma.

**Lemma 9.7.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type I' vicinity, we have*

$$F(\hat{X}, \hat{B}, V, T) = O\left(\frac{1}{(N\eta)^{n+2}}\right). \quad (9.35)$$

Recalling the functions  $\mathbb{G}(\hat{B}, T)$  and  $\mathbb{F}(\hat{X}, V)$  defined in (6.23) and (6.24), we further introduce

$$\mathring{\mathbb{G}}(\hat{B}, T) = \exp\{(a_+ - a_-)N\eta\}\mathbb{G}(\hat{B}, T), \quad \mathring{\mathbb{F}}(\hat{X}, V) = \exp\{-(a_+ - a_-)N\eta\}\mathbb{F}(\hat{X}, V). \quad (9.36)$$

Then, we have the decomposition

$$F(\hat{X}, \hat{B}, V, T) = \mathring{\mathbb{G}}(\hat{B}, T)\mathring{\mathbb{F}}(\hat{X}, V). \quad (9.37)$$

Hence, we can estimate  $\mathring{\mathbb{F}}(\hat{X}, V)$  and  $\mathring{\mathbb{G}}(\hat{B}, T)$  separately in the sequel.

9.2.1. *Estimate of  $\mathring{\mathbb{F}}(\hat{X}, V)$ .* We have the following lemma.

**Lemma 9.8.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type I' vicinity, we have*

$$\mathring{\mathbb{F}}(\hat{X}, V) = O\left(\frac{1}{N\eta}\right). \quad (9.38)$$

*Proof.* Using (8.8) and (8.9), we can write

$$X_j = P_1^* V_j^* \hat{X}_j V_j P_1 = P_1^* D_{\pm} P_1 + O\left(\frac{\Theta}{\sqrt{M}}\right), \quad (9.39)$$

where the remainder term represents a  $2 \times 2$  matrix whose max-norm is bounded by  $\Theta/\sqrt{M}$ . Using (9.39) and recalling  $N = MW$  yields

$$\exp\left\{M\eta \sum_{j=1}^W \text{Tr} X_j J\right\} = \exp\left\{N\eta \text{Tr} P_1^* D_{\pm} P_1 J\right\} \left(1 + O\left(\frac{\Theta N\eta}{\sqrt{M}}\right)\right). \quad (9.40)$$

Substituting (9.40) into (3.19) and (6.24), we can write

$$\begin{aligned} \mathring{\mathbb{F}}(\hat{X}, V) &= \int d\mu(P_1) dX^{[1]} \exp\left\{N\eta \text{Tr} P_1^* D_{\pm} P_1 J\right\} \\ &\quad \times \prod_{k=p,q} \frac{1}{\det^2(X_k^{[1]})} \cdot \exp\left\{i \text{Tr} X_k^{[1]} J Z - \sum_j \tilde{\mathfrak{s}}_{jk} \text{Tr} X_j X_k^{[1]} J\right\} \\ &\quad \times \prod_{k,\ell=p,q} \exp\left\{\frac{\tilde{\mathfrak{s}}_{k\ell}}{2M} \text{Tr} X_k^{[1]} J X_{\ell}^{[1]} J\right\} \cdot \left(1 + O\left(\frac{\Theta N\eta}{\sqrt{M}}\right)\right). \end{aligned}$$

Recalling the parametrization of  $P_1$  in (3.24), we have

$$\text{Tr} P_1^* D_{\pm} P_1 J = (1 - 2v^2)(a_+ - a_-).$$

Consequently, we have

$$\begin{aligned} \mathring{\mathbb{F}}(\hat{X}, V) &= \int dX^{[1]} \int v dv \int \frac{d\theta}{\pi} \exp\left\{-2(a_+ - a_-)N\eta v^2\right\} \prod_{k=p,q} \frac{1}{\det^2(X_k^{[1]})} \\ &\quad \times \prod_{k=p,q} \exp\left\{i \text{Tr} X_k^{[1]} J Z - \sum_j \tilde{\mathfrak{s}}_{jk} \text{Tr} X_j X_k^{[1]} J\right\} (1 + o(1)). \end{aligned}$$

Obviously, by the fact that  $X^{[1]}$ -variables are all bounded and  $|\det X_k^{[1]}| = 1$  for  $k = p, q$ , it is easy to see that

$$|\mathring{\mathbb{F}}(\hat{X}, V)| \leq C \int_0^1 v dv \exp \left\{ -2(a_+ - a_-)N\eta v^2 \right\} = O\left(\frac{1}{N\eta}\right).$$

Therefore, we completed the proof.  $\square$

9.2.2. *Estimate of  $\mathring{\mathbb{G}}(\hat{B}, T)$ .* Recall the definition of  $\mathring{\mathbb{G}}(\hat{B}, T)$  from (9.36), (6.23) and (3.20). In this section, we will prove the following lemma.

**Lemma 9.9.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type I' vicinity, we have*

$$\mathring{\mathbb{G}}(\hat{B}, T) = O\left(\frac{1}{(N\eta)^{n+1}}\right). \quad (9.41)$$

Note that  $y_p^{[1]}, y_q^{[1]}$  and  $t$  in the parametrization of  $Q_1$  (see (3.24)) are not bounded, we shall truncate them with some appropriate bounds at first, whereby we can neglect some irrelevant terms in the integrand, in order to simplify the integral. More specifically, we will do the truncations

$$t \leq \frac{1}{(N\eta)^{1/4}} \quad (9.42)$$

and

$$y_p^{[1]}, y_q^{[1]} \leq (N\eta)^{\frac{1}{8}}. \quad (9.43)$$

Accordingly, we set

$$\begin{aligned} \widehat{\mathbb{G}}(\hat{B}, T) &:= e^{(a_+ - a_-)N\eta} \int_{\mathbb{L}} \frac{d\sigma}{2\pi} \int_{\mathbb{I}^2} v_p^{[1]} v_q^{[1]} dv_p^{[1]} dv_q^{[1]} \int_0^{(N\eta)^{\frac{1}{8}}} dy_p^{[1]} \int_0^{(N\eta)^{\frac{1}{8}}} dy_q^{[1]} \\ &\quad \times \int_0^{(N\eta)^{-\frac{1}{4}}} 2t dt \int_{\mathbb{L}^2} d\sigma_p^{[1]} d\sigma_q^{[1]} g(Q_1, T, \hat{B}, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}), \end{aligned} \quad (9.44)$$

where we have used the parameterization of  $\mathbf{w}^{[1]}$  in (3.15). We will prove the following lemma.

**Lemma 9.10.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type I' vicinity, we have*

$$\mathring{\mathbb{G}}(\hat{B}, T) = \widehat{\mathbb{G}}(\hat{B}, T) + O(e^{-N^\varepsilon})$$

for some positive constant  $\varepsilon$ .

*Proof.* At first, by (6.29)-(6.32), we have for any  $j$ ,

$$\text{Re} Tr B_j Y_k^{[1]} J \geq \frac{y_k^{[1]}}{(s+t)^2} \cdot \frac{\min\{\text{Re} b_{j,1}, \text{Re} b_{j,2}\}}{(s_j + t_j)^2} \geq c \frac{y_k^{[1]}}{1 + 2t^2}, \quad k = p, q, \quad (9.45)$$

for some positive constant  $c$ , where the last step follows from the facts that  $\text{Re} b_{j,1}, \text{Re} b_{j,2} = \text{Re} a_+ + o(1)$  and  $t_j = o(1)$  in the Type I' vicinity. In addition, it is not difficult to get

$$Tr B_j J = \left( a_+ - a_- + O\left(\frac{\Theta}{\sqrt{M}}\right) \right) (1 + 2t^2), \quad \forall j = 1, \dots, W,$$

which implies that

$$M\eta \sum_{j=1}^W Tr B_j J = (a_+ - a_-)N\eta + 2\left(a_+ - a_- + O\left(\frac{\Theta}{\sqrt{M}}\right)\right)N\eta t^2 + O\left(\frac{\Theta N\eta}{\sqrt{M}}\right). \quad (9.46)$$

Note that the second and third factors in the definition of  $g(\cdot)$  in (3.20) can be bounded by 1, according to (6.26). Then, as a consequence of (9.45) and (9.46), we have

$$e^{(a_+ - a_-)N\eta} |g(\cdot)| \leq C(y_p^{[1]} y_q^{[1]})^{n+3} \exp\{-c'N\eta t^2\} \exp\left\{-c \frac{y_p^{[1]} + y_q^{[1]}}{1 + 2t^2}\right\}, \quad (9.47)$$

for some positive constants  $C$ ,  $c$  and  $c'$ . By integrating  $y_p^{[1]}$  and  $y_q^{[1]}$  out at first, we can easily see that the truncation (9.42) only produces an error of order  $O(\exp\{-N^\varepsilon\})$  to the integral  $\mathring{\mathbb{G}}(\hat{B}, T)$ , for some positive constant  $\varepsilon = \varepsilon(\varepsilon_2)$  by the assumption  $\eta \geq N^{-1+\varepsilon_2}$  in (1.15). Then one can substitute the bound (9.42) to the last factor of the r.h.s. of (9.47), thus

$$\exp\left\{-c\frac{y_p^{[1]} + y_q^{[1]}}{1 + 2t^2}\right\} \leq \exp\left\{-\frac{c}{2}(y_p^{[1]} + y_q^{[1]})\right\}.$$

We can also do the truncation (9.43) in the integral  $\mathring{\mathbb{G}}(\hat{B}, T)$ , up to an error of order  $O(\exp\{-N^\varepsilon\})$ , for some positive constant  $\varepsilon$ . Therefore, we completed the proof of Lemma 9.10.  $\square$

With the aid of Lemma 9.10, it suffices to work on  $\widehat{\mathbb{G}}(\hat{B}, T)$  in the sequel. We have the following lemma.

**Lemma 9.11.** *We have*

$$\widehat{\mathbb{G}}(\hat{B}, T) = O\left(\frac{1}{(N\eta)^{n+1}}\right).$$

*Proof of Lemma 9.11.* Recall the parameterization of  $\mathbf{w}_k^{[1]}$  in (3.15) again. To simplify the notation, we set

$$\mathbf{w}_k^{[1]} = u_k^{[1]}v_k^{[1]}, \quad k = p, q.$$

Similarly to (9.39), using  $t = o(1)$  from (9.42), we have the expansion

$$B_j = Q_1^{-1}T_j^{-1}\hat{B}_jT_jQ_1 = Q_1^{-1}D_\pm Q_1 + O\left(\frac{\Theta}{\sqrt{M}}\right).$$

Consequently, we have

$$-M\eta \sum_{j=1}^W \text{Tr} B_j J = -N\eta(a_+ - a_-)(1 + 2t^2) + O\left(\frac{\Theta N\eta}{\sqrt{M}}\right). \quad (9.48)$$

In addition, for  $k = p, q$ , using the fact  $\sum_j \tilde{\mathbf{s}}_{jk} = 1$ , we have

$$\begin{aligned} \sum_j \tilde{\mathbf{s}}_{jk} \text{Tr} B_j Y_k^{[1]} J &= \text{Tr} Q_1^{-1} D_\pm Q_1 Y_k^{[1]} J + \frac{\Theta}{\sqrt{M}} \text{Tr} R_k Y_k^{[1]} \\ &= y_k^{[1]} \left( (a_+ - a_-)t^2 + (a_+(u_k^{[1]})^2 - a_-(v_k^{[1]})^2) \right) \\ &\quad + y_k^{[1]} \left( (a_+ - a_-)(e^{-i(\sigma_k^{[1]} + \sigma)} + e^{i(\sigma_k^{[1]} + \sigma)}) \mathbf{w}_k^{[1]} st \right) + \frac{\Theta}{\sqrt{M}} \text{Tr} R_k Y_k^{[1]}, \end{aligned} \quad (9.49)$$

where  $R_k$  is a  $2 \times 2$  matrix independent of  $Y_k^{[1]}$ , satisfying

$$\|R_k\|_{\max} = O(1). \quad (9.50)$$

Observe that the term in (9.48) is obviously independent of  $\mathbf{w}^{[1]}$ -variables. In addition, for  $k = p$  or  $q$ , we have

$$i \text{Tr} Y_k^{[1]} J Z = (-\eta + iE(1 - 2(v_k^{[1]})^2))y_k^{[1]}, \quad (9.51)$$

and for  $k, \ell = p$  or  $q$ , we have

$$\text{Tr} Y_k^{[1]} J Y_\ell^{[1]} J = y_k^{[1]} y_\ell^{[1]} \left( (\mathbf{w}_k^{[1]} \mathbf{w}_\ell^{[1]})^2 + \mathbf{w}_k^{[1]} \mathbf{w}_\ell^{[1]} \left( e^{i(\sigma_k^{[1]} - \sigma_\ell^{[1]})} - e^{i(\sigma_\ell^{[1]} - \sigma_k^{[1]})} \right) \right). \quad (9.52)$$

Moreover, we have

$$\left( \left( \mathbf{w}_q^{[1]} (\mathbf{w}_q^{[1]})^* \right)_{12} \left( \mathbf{w}_p^{[1]} (\mathbf{w}_p^{[1]})^* \right)_{21} \right)^n = \left( \mathbf{w}_p^{[1]} \mathbf{w}_q^{[1]} \right)^n e^{in(\sigma_p^{[1]} - \sigma_q^{[1]})}. \quad (9.53)$$

Substituting (9.48), (9.49) and (9.51)-(9.53) to the definition of  $g(\cdot)$  in (3.20) and reordering the factors properly, we can write the integrand in (9.44) as

$$\begin{aligned}
& \exp\{(a_+ - a_-)N\eta\}g(\cdot) \\
&= \exp\{in(\sigma_p^{[1]} - \sigma_q^{[1]})\} \cdot \exp\left\{- (a_+ - a_-)st \sum_{k=p,q} y_k^{[1]} \mathbf{w}_k^{[1]} \left(e^{-i(\sigma_k^{[1]} + \sigma)} + e^{i(\sigma_k^{[1]} + \sigma)}\right)\right\} \\
&\times \exp\left\{-\frac{\Theta}{\sqrt{M}} \sum_{k=p,q} Tr R_k Y_k^{[1]}\right\} \cdot \exp\left\{-\frac{1}{M} \tilde{s}_{pq} y_p^{[1]} y_q^{[1]} \mathbf{w}_p^{[1]} \mathbf{w}_q^{[1]} \left(e^{i(\sigma_p^{[1]} - \sigma_q^{[1]})} - e^{i(\sigma_q^{[1]} - \sigma_p^{[1]})}\right)\right\} \\
&\times \prod_{k=p,q} (y_k^{[1]})^{n+3} (\mathbf{w}_k^{[1]})^n \cdot \prod_{k,\ell=p,q} \exp\left\{-\frac{1}{2M} \tilde{s}_{k\ell} y_k^{[1]} y_\ell^{[1]} \left(\mathbf{w}_k^{[1]} \mathbf{w}_\ell^{[1]}\right)^2\right\} \cdot \exp\{-2N\eta(a_+ - a_-)t^2\} \\
&\times \prod_{k=p,q} \exp\left\{-y_k^{[1]} \left((a_+(u_k^{[1]})^2 - a_-(v_k^{[1]})^2) + (a_+ - a_-)t^2 + \eta - iE(1 - 2(v_k^{[1]})^2)\right)\right\} \\
&\times \left(1 + O\left(\frac{\Theta N \eta}{\sqrt{M}}\right)\right),
\end{aligned} \tag{9.54}$$

where the last factor is independent of the  $\mathbf{w}^{[1]}$ -variables. Here, we put the factors containing  $\sigma_p^{[1]}$  and  $\sigma_q^{[1]}$  together, namely, the first two lines on the r.h.s. of (9.54).

For further discussion, we write for  $k = p, q$

$$Tr R_k Y_k^{[1]} = y_k^{[1]} (\mathbf{r}_k^+ e^{i\sigma_k} + \mathbf{r}_k^- e^{-i\sigma_k} + \mathbf{r}_k), \tag{9.55}$$

where  $\mathbf{r}_k^+$ ,  $\mathbf{r}_k^-$  and  $\mathbf{r}_k$  are all polynomials of  $u_k^{[1]}$  and  $v_k^{[1]}$ , with bounded degree and bounded coefficients, in light of (9.50), the definition of  $Y_k^{[1]}$  in (3.14) and the parametrization in (3.15).

Now, we start to estimate the integral (9.44) by using (9.54). We deal with the integral over  $\sigma_p^{[1]}$  and  $\sigma_q^{[1]}$  at first. These variables are collected in the integral of the form

$$\begin{aligned}
\mathcal{I}_\sigma(\ell_1, \ell_2) &:= \int_{\mathbb{L}^2} d\sigma_p^{[1]} d\sigma_q^{[1]} \exp\{i(n + \ell_1)\sigma_p^{[1]}\} \exp\{-i(n + \ell_2)\sigma_q^{[1]}\} \exp\left\{-\frac{\Theta}{\sqrt{M}} \sum_{k=p,q} Tr R_k Y_k^{[1]}\right\} \\
&\times \exp\left\{-(a_+ - a_-)st \sum_{k=p,q} y_k^{[1]} \mathbf{w}_k^{[1]} (e^{-i(\sigma_k^{[1]} + \sigma)} + e^{i(\sigma_k^{[1]} + \sigma)})\right\} \\
&\times \exp\left\{-\frac{1}{M} \tilde{s}_{pq} y_p^{[1]} y_q^{[1]} \mathbf{w}_p^{[1]} \mathbf{w}_q^{[1]} (e^{i(\sigma_p^{[1]} - \sigma_q^{[1]})} - e^{i(\sigma_q^{[1]} - \sigma_p^{[1]})})\right\}
\end{aligned}$$

with integers  $\ell_1$  and  $\ell_2$  independent of  $n$ . Note that according to (9.54), it suffices to consider  $\mathcal{I}_\sigma(0, 0)$  for the proof of (9.41). We study  $\mathcal{I}_\sigma(\ell_1, \ell_2)$  for general  $\ell_1$  and  $\ell_2$  here, which will be used later.

Now, we set

$$\begin{aligned}
c_{p,q} &:= \tilde{s}_{pq} y_p^{[1]} y_q^{[1]} \mathbf{w}_p^{[1]} \mathbf{w}_q^{[1]}, \\
c_{k,1} &:= -(a_+ - a_-)st y_k^{[1]} \mathbf{w}_k^{[1]} e^{-i\sigma} - \frac{\Theta}{\sqrt{M}} y_k^{[1]} \mathbf{r}_k^-, \quad k = p, q, \\
c_{k,2} &:= -(a_+ - a_-)st y_k^{[1]} \mathbf{w}_k^{[1]} e^{i\sigma} - \frac{\Theta}{\sqrt{M}} y_k^{[1]} \mathbf{r}_k^+, \quad k = p, q.
\end{aligned} \tag{9.56}$$

In addition, we introduce

$$d_{p,q} := y_p^{[1]} y_q^{[1]}, \quad d_k := \left(t + \frac{\Theta}{\sqrt{M}}\right) y_k^{[1]}, \quad k = p, q. \tag{9.57}$$

Obviously, when (9.42) is satisfied, we have

$$c_{p,q} = O(d_{p,q}), \quad c_{k,1} = O(d_k), \quad c_{k,2} = O(d_k), \quad k = p, q. \tag{9.58}$$

With the aid of the notation defined in (9.55) and (9.56), we can write

$$\begin{aligned} \mathcal{I}_\sigma(\ell_1, \ell_2) &= \exp \left\{ -\frac{\Theta}{\sqrt{M}}(y_p^{[1]}\mathbf{r}_p + y_q^{[1]}\mathbf{r}_q) \right\} \int_{\mathbb{L}^2} d\sigma_p d\sigma_q \exp \left\{ \mathbf{i}(n + \ell_1)\sigma_p^{[1]} \right\} \exp \left\{ -\mathbf{i}(n + \ell_2)\sigma_q^{[1]} \right\} \\ &\quad \times \prod_{k=p,q} \exp \left\{ c_{k,1}e^{-\mathbf{i}\sigma_k^{[1]}} + c_{k,2}e^{\mathbf{i}\sigma_k^{[1]}} \right\} \exp \left\{ -\frac{c_{p,q}}{M}e^{\mathbf{i}(\sigma_p^{[1]} - \sigma_q^{[1]})} + \frac{c_{p,q}}{M}e^{\mathbf{i}(\sigma_q^{[1]} - \sigma_p^{[1]})} \right\}. \end{aligned} \quad (9.59)$$

We have the following lemma.

**Lemma 9.12.** *Under the truncation (9.42) and (9.43), we have*

$$|\mathcal{I}_\sigma(\ell_1, \ell_2)| \leq C \left( \left( \frac{d_{p,q}}{M} \right)^{n+\ell_3} + d_p^{2(n+\ell_3)} + d_q^{2(n+\ell_3)} \right), \quad \ell_3 := \frac{\ell_1 + \ell_2}{2}$$

for some positive constant  $C$ .

*Proof.* At first, by Taylor expansion, we have

$$\begin{aligned} &\exp \left\{ \mathbf{i}(n + \ell_1)\sigma_p^{[1]} \right\} \exp \left\{ -\mathbf{i}(n + \ell_2)\sigma_q^{[1]} \right\} \exp \left\{ -\frac{c_{p,q}}{M}e^{\mathbf{i}(\sigma_p^{[1]} - \sigma_q^{[1]})} + \frac{c_{p,q}}{M}e^{\mathbf{i}(\sigma_q^{[1]} - \sigma_p^{[1]})} \right\} \\ &= \sum_{n_1, n_2=0}^{\infty} \frac{(-1)^{n_1}}{(n_1)!(n_2)!} \left( \frac{c_{p,q}}{M} \right)^{n_1+n_2} \exp \left\{ \mathbf{i}(n + \ell_1 + n_1 - n_2)\sigma_p^{[1]} \right\} \exp \left\{ -\mathbf{i}(n + \ell_2 + n_1 - n_2)\sigma_q^{[1]} \right\}. \end{aligned} \quad (9.60)$$

Now, for any  $m_1, m_2 \in \mathbb{Z}$ , we denote

$$\begin{aligned} \tilde{\mathcal{I}}_\sigma(m_1, m_2) &:= \int_{\mathbb{L}^2} d\sigma_p^{[1]} d\sigma_q^{[1]} \exp \{ \mathbf{i}m_1\sigma_p^{[1]} \} \exp \{ -\mathbf{i}m_2\sigma_q^{[1]} \} \prod_{k=p,q} \exp \left\{ c_{k,1}e^{-\mathbf{i}\sigma_k^{[1]}} + c_{k,2}e^{\mathbf{i}\sigma_k^{[1]}} \right\} \\ &= 4\pi^2 \sum_{n_3=0}^{\infty} \mathbf{1}(n_3 + m_1 \geq 0) \frac{(c_{p,1})^{n_3+m_1} (c_{p,2})^{n_3}}{n_3!(n_3 + m_1)!} \sum_{n_4=0}^{\infty} \mathbf{1}(n_4 + m_2 \geq 0) \frac{(c_{q,1})^{n_4} (c_{q,2})^{n_4+m_2}}{n_4!(n_4 + m_2)!}. \end{aligned} \quad (9.61)$$

Setting

$$m_1 := n + \ell_1 + n_1 - n_2, \quad m_2 := n + \ell_2 + n_1 - n_2, \quad (9.62)$$

and using (9.60), we can rewrite (9.59) as

$$\mathcal{I}_\sigma(\ell_1, \ell_2) = \exp \left\{ -\frac{\Theta}{\sqrt{M}}(y_p^{[1]}\mathbf{r}_p + y_q^{[1]}\mathbf{r}_q) \right\} \sum_{n_1, n_2=0}^{\infty} \frac{(-1)^{n_1}}{(n_1)!(n_2)!} \left( \frac{c_{p,q}}{M} \right)^{n_1+n_2} \tilde{\mathcal{I}}_\sigma(m_1, m_2). \quad (9.63)$$

For simplicity, we employ the notation

$$m_3 := m_3(\ell_1, n_1, n_2, n_3) = m_1 + n_3, \quad m_4 := m_4(\ell_2, n_1, n_2, n_4) = m_2 + n_4. \quad (9.64)$$

Consequently, by (9.63) and (9.61) we obtain

$$\begin{aligned} |\mathcal{I}_\sigma(\ell_1, \ell_2)| &\leq 4\pi^2 \left| \exp \left\{ -\frac{\Theta}{\sqrt{M}}(y_p^{[1]}\mathbf{r}_p + y_q^{[1]}\mathbf{r}_q) \right\} \right| \sum_{n_1, n_2=0}^{\infty} \frac{1}{(n_1)!(n_2)!} \left| \frac{c_{p,q}}{M} \right|^{n_1+n_2} \\ &\quad \times \sum_{n_3=0}^{\infty} \mathbf{1}(m_3 \geq 0) \frac{|c_{p,1}|^{m_3} |c_{p,2}|^{n_3}}{n_3!m_3!} \cdot \sum_{n_4=0}^{\infty} \mathbf{1}(m_4 \geq 0) \frac{|c_{q,1}|^{n_4} |c_{q,2}|^{m_4}}{n_4!m_4!} \\ &\leq C \max_{n_1, n_2, n_3, n_4, m_3, m_4 \in \mathbb{N}} \left| \left( \frac{c_{p,q}}{M} \right)^{n_1+n_2} (c_{p,1})^{m_3} (c_{p,2})^{n_3} (c_{q,1})^{n_4} (c_{q,2})^{m_4} \right| \\ &\leq C \max_{n_1, n_2, n_3, n_4, m_3, m_4 \in \mathbb{N}} \left| \left( \frac{c_{p,q}}{M} \right)^{n_1+n_2} (c_{p,1})^{m_3} (c_{q,2})^{m_4} \right| \end{aligned} \quad (9.65)$$

for some positive constant  $C$ , where in the last step we used the fact

$$|c_{k,1}| < 1, \quad |c_{k,2}| < 1, \quad k = p, q \quad (9.66)$$

which can be seen directly from the definition in (9.58), the truncations (9.42) and (9.43) and the assumption  $\eta \leq M^{-1}N^{\varepsilon_2}$ . Analogously, we also have

$$\left| \frac{c_{p,q}}{M} \right| < 1. \quad (9.67)$$

According to the definitions (9.62) and (9.64), we have

$$2(n_1 + n_2) + m_3 + m_4 \geq 2n + \ell_1 + \ell_2.$$

Hence, by using (9.66) and (9.67), we have the trivial bound

$$\max_{n_1, n_2, n_3, n_4, m_3, m_4 \geq 0} \left| \left( \sqrt{\frac{|c_{p,q}|}{M}} \right)^{2(n_1+n_2)} (c_{p,1})^{m_3} (c_{q,2})^{m_4} \right| \leq \left( \frac{|c_{p,q}|}{M} \right)^{n+\ell_3} + |c_{p,1}|^{2(n+\ell_3)} + |c_{q,2}|^{2(n+\ell_3)}.$$

Therefore, we completed the proof by using (9.58).  $\square$

Now, we return to the proof of Lemma 9.11. Using (9.54) and Lemma 9.12 with  $\ell_1 = \ell_2 = 0$  to (9.44), and integrating the bounded variables  $v_p^{[1]}$ ,  $v_q^{[1]}$  and  $\sigma$  out, we can get

$$\begin{aligned} |\widehat{\mathbb{G}}(\hat{B}, T)| &\leq C \int_0^{(N\eta)^{\frac{1}{8}}} dy_p^{[1]} \int_0^{(N\eta)^{\frac{1}{8}}} dy_q^{[1]} \int_0^{(N\eta)^{-\frac{1}{4}}} 2t dt \cdot (d_{p,q})^{n+3} \\ &\quad \times \left( \left( \frac{d_{p,q}}{M} \right)^n + d_p^{2n} + d_q^{2n} \right) \exp \{ -2N\eta(a_+ - a_-)t^2 \} \\ &\quad \times \exp \left\{ -\frac{\sqrt{4-E^2}}{2} \sum_{k=p,q} y_k^{[1]} \right\} (1 + o(1)) \end{aligned}$$

where the last two factors come from the facts

$$\begin{aligned} \left| \exp \left\{ -\sum_{k=p,q} y_k^{[1]} \left( a_+(u_k^{[1]})^2 - a_-(v_k^{[1]})^2 \right) \right\} \right| &= \exp \left\{ -\frac{\sqrt{4-E^2}}{2} \sum_{k=p,q} y_k^{[1]} \right\}, \\ \left| \exp \left\{ -\sum_{k=p,q} y_k^{[1]} \left( (a_+ - a_-)t^2 + \eta - \mathbf{i}E(1 - 2(v_k^{[1]})^2) \right) \right\} \right| &= 1 + o(1). \end{aligned} \quad (9.68)$$

In (9.68) we used the fact  $(u_k^{[1]})^2 + (v_k^{[1]})^2 = 1$ . Now, we integrate  $y_p^{[1]}$  and  $y_q^{[1]}$  out. Consequently, by the definition in (9.57), we have

$$|\widehat{\mathbb{G}}(\hat{B}, T)| \leq C \int_0^{(N\eta)^{-\frac{1}{4}}} 2t dt \left( \frac{1}{M^n} + \left( \frac{\Theta}{\sqrt{M}} \right)^{2n} + t^{2n} \right) \exp \{ -2N\eta(a_+ - a_-)t^2 \} = O\left( \frac{1}{(N\eta)^{n+1}} \right),$$

where in the last step we have used the assumption  $\eta \leq M^{-1}N^{\varepsilon_2}$  in (1.15), Assumption 1.13, the definition of  $\Theta$  in (5.33) and the fact  $N = MW$ . Hence, we completed the proof of Lemma 9.11.  $\square$

Finally, we can prove Lemma 9.9, and further prove Lemma 9.7.

*Proof of Lemma 9.9.* This is a direct consequence of Lemmas 9.10 and 9.11.  $\square$

*Proof of Lemma 9.7.* This is a direct consequence of (9.37), Lemma 9.8 and Lemma 9.9.  $\square$

**9.3. Summing up: Proof of Lemma 9.1.** In this section, we slightly modify the discussions in Sections 9.1 and 9.2 to prove Lemma 9.1. The combination of Lemmas 9.2 and 9.7 would directly imply Lemma 9.1 if the  $\mathbf{Q}(\cdot)$  factor were not present in the definition of  $\mathbf{A}(\cdot)$ . Now we should take  $\mathbf{Q}(\cdot)$  into account. This argument is similar to the corresponding discussion in Section 6.4.

*Proof of Lemma 9.1.* At first, we observe that  $\kappa_1$ ,  $\kappa_2$  and  $\kappa_3$  in (6.5) are obviously independent of  $n$ . Then, by the fact  $\kappa_1 = W^{O(1)}$ , it suffices to consider one monomial of the form

$$\mathfrak{p}_1\left(t, s, (y_p^{[1]})^{-1}, (y_q^{[1]})^{-1}\right) \mathfrak{p}_2\left(\left\{e^{i\sigma_k^{[1]}}, e^{-i\sigma_k^{[1]}}\right\}_{k=p,q}\right) \mathfrak{q}\left(\left\{\frac{\omega_{i,a}\xi_{j,b}}{M}\right\}_{\substack{i,j=1,\dots,W \\ a,b=1,\dots,4}}\right),$$

where the degrees of  $\mathfrak{p}_1(\cdot)$ ,  $\mathfrak{p}_2(\cdot)$  and  $\mathfrak{q}(\cdot)$  are all  $O(1)$ , and independent of  $n$ , in light of the fact  $\kappa_3 = O(1)$  in (6.5). Especially, the order of  $(y_p^{[1]})^{-1}$  and  $(y_q^{[1]})^{-1}$  are not larger than 2, which can be easily seen from the definition of  $\mathcal{Q}(\cdot)$  in (3.18).

Now, we reuse the notation  $\mathsf{P}_q(\hat{X}, \hat{B}, V, T)$  and  $\mathsf{F}_p(\hat{X}, \hat{B}, V, T)$  in (6.40), by redefining them as

$$\begin{aligned} \mathsf{P}_q(\hat{X}, \hat{B}, V, T) &:= \int d\Omega d\Xi \mathcal{P}(\cdot) \cdot \mathfrak{q}\left(\left\{\frac{\omega_{i,a}\xi_{j,b}}{M}\right\}_{\substack{i,j=1,\dots,W \\ a,b=1,\dots,4}}\right), \\ \mathsf{F}_p(\hat{X}, \hat{B}, V, T) &:= \int dX^{[1]} d\mathbf{Y}^{[1]} d\mathbf{W}^{[1]} d\mu(P_1) d\nu(Q_1) \mathcal{F}(\cdot) \\ &\quad \times \mathfrak{p}_1\left(t, s, (y_p^{[1]})^{-1}, (y_q^{[1]})^{-1}\right) \mathfrak{p}_2\left(\left\{e^{i\sigma_k^{[1]}}, e^{-i\sigma_k^{[1]}}\right\}_{k=p,q}\right). \end{aligned}$$

It is easy to check  $\mathcal{P}(\cdot)\mathfrak{q}(\cdot)$  also has an expansion of the form in (9.14). Hence,  $\mathsf{P}_q(\cdot)$  can be bounded in the same way as  $\mathcal{P}(\cdot)$ , thus we have

$$|\mathsf{P}_q(\hat{X}, \hat{B}, V, T)| \leq \frac{W^{2+\gamma}\Theta^2}{M} |\det \mathbb{A}_+|^2 \det(S^{(1)})^2.$$

For  $\mathsf{F}_p(\cdot)$ , the main modification is to use Lemma 9.12 with general  $\ell_1$  and  $\ell_2$  independent of  $n$ , owing to the function  $\mathfrak{p}_2(\cdot)$ . In addition, by the truncations in (9.42) and (9.43), we can bound  $\mathfrak{p}_1(\cdot)$  by some constant  $C$ . Hence, it suffices to replace  $n$  by  $n + \ell_3$  in the proof of Lemma 9.11. Finally, we can get

$$\mathsf{F}_p(\hat{X}, \hat{B}, V, T) = O\left(\frac{1}{(N\eta)^{n+\ell_3}}\right),$$

with some finite integer  $\ell_3$  independent of  $n$ . Consequently, we completed the proof of Lemma 9.1.  $\square$

## 10. INTEGRAL OVER THE TYPE II AND III VICINITIES

In this section, we prove Lemma 5.9. We only present the discussion for  $\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \mathbb{I}^{W-1})$ , i.e. integral over the Type II vicinity. The discussion on  $\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_-^x, \Upsilon_S, \mathbb{I}^{W-1})$  is analogous. We start from (8.48). In this section, we will regard  $V$ -variables as fixed parameters, and consider the integrand as a function of all the other variables.

Similarly, we shall provide an estimate for the integrand. At first, under the parameterization (8.2) with  $\varkappa = +$ , we see that

$$x_{j,1} - x_{j,2} = \frac{ia_+}{\sqrt{M}}(\hat{x}_{j,1} - \hat{x}_{j,2}) + O\left(\frac{\Theta}{M}\right), \quad b_{j,1} + b_{j,2} = a_+ - a_- + O\left(\frac{\Theta^{\frac{1}{2}}}{\sqrt{M}}\right).$$

Consequently, we have

$$\begin{aligned} \prod_{j=1}^W (x_{j,1} - x_{j,2})^2 (b_{j,1} + b_{j,2})^2 &= \frac{(-a_+^2)^W}{M^W} (a_+ - a_-)^{2W} \left(1 + O\left(\frac{\Theta^{\frac{3}{2}}}{\sqrt{M}}\right)\right) \\ &\quad \times \prod_{j=1}^W \left(\hat{x}_{j,1} - \hat{x}_{j,2} + O\left(\frac{\Theta}{\sqrt{M}}\right)\right)^2. \end{aligned} \quad (10.1)$$

Now, analogously to the case of Type I vicinity, what remains is to estimate  $\mathsf{A}(\hat{X}, \hat{B}, V, T)$ . Our aim is to prove the following lemma.

**Lemma 10.1.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type II vicinity, we have*

$$|A(\hat{X}, \hat{B}, V, T)| \leq e^{-cN\eta} |\det \mathbb{A}_+|^2 \det(S^{(1)})^2 \quad (10.2)$$

for some positive constant  $c$ .

With the aid of (10.1) and Lemma 10.1, we can prove Lemma 5.9.

*Proof of Lemma 5.9.* Recall (8.48). At first, by the definition of  $\mathbb{A}_+^v$  in (8.35), (5.23) and the fact  $\operatorname{Re} a_+^2 > 0$ , we can see that

$$\operatorname{Re}(\mathring{\mathbf{x}}' \mathbb{A}_+^v \mathring{\mathbf{x}}) \geq \|\mathring{\mathbf{x}}\|_2^2 \quad (10.3)$$

for all  $\{V_j\}_{j=2}^W \in (\hat{U}(2))^{W-1}$ . Substituting (5.20), (10.1), (10.2), (10.3) and the estimates in Proposition 8.6 into (8.48) yields

$$\begin{aligned} & |\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_+^x, \Upsilon_S, \mathbb{I}^{W-1})| \\ & \leq e^{-cN\eta} \cdot \frac{(a_+ - a_-)^{2W}}{8^W \pi^{3W-1}} \cdot |\det S^{(1)}|^2 \cdot |\det \mathbb{A}_+| \cdot \int_{\mathbb{L}^{W-1}} \prod_{j=2}^W \frac{d\theta_j}{2\pi} \int_{\mathbb{I}^{W-1}} \prod_{j=2}^W 2v_j dv_j \\ & \times \int_{\mathbb{R}^{W-1}} \prod_{j=2}^W d\tau_{j,1} \int_{\mathbb{R}^{W-1}} \prod_{j=2}^W d\tau_{j,2} \int_{\mathbb{R}^W} \prod_{j=1}^W dc_{j,1} \int_{\mathbb{R}^W} \prod_{j=1}^W dc_{j,2} \int_{\mathbb{R}^W} \prod_{j=1}^W d\hat{x}_{j,1} \int_{\mathbb{R}^W} \prod_{j=1}^W d\hat{x}_{j,2} \\ & \times \exp\{(a_+ - a_-)^2 \boldsymbol{\tau}'_1 S^{(1)} \boldsymbol{\tau}_1\} \exp\{(a_+ - a_-)^2 \boldsymbol{\tau}'_2 S^{(1)} \boldsymbol{\tau}_2\} \exp\{-\frac{1}{2} \|\mathring{\mathbf{c}}_1\|_2^2 - \frac{1}{2} \|\mathring{\mathbf{c}}_2\|_2^2\} \\ & \times \exp\{-\frac{1}{2} \|\mathring{\mathbf{x}}_1\|_2^2 - \frac{1}{2} \|\mathring{\mathbf{x}}_2\|_2^2\} \prod_{j=1}^W \left( \hat{x}_{j,1} - \hat{x}_{j,2} + O\left(\frac{\Theta}{\sqrt{M}}\right) \right)^2, \end{aligned}$$

where we absorbed several factors by  $\exp\{-cN\eta\}$ . We also enlarged the domains to the full ones. Then, using the trivial facts

$$\int_{\mathbb{L}^{W-1}} \prod_{j=2}^W \frac{d\theta_j}{2\pi} \int_{\mathbb{I}^{W-1}} \prod_{j=2}^W 2v_j dv_j = 1$$

and performing the Gaussian integral for the remaining variables, we can get

$$|\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_+^x, \Upsilon_S, \mathbb{I}^{W-1})| \leq C |\det S^{(1)}| \cdot |\det \mathbb{A}_+| \cdot \left(1 + O\left(\frac{\Theta}{\sqrt{M}}\right)\right)^W. \quad (10.4)$$

Observe that

$$|\det \mathbb{A}_+| \leq |1 + a_+^2|^W \leq 2^W. \quad (10.5)$$

Moreover, by Assumption 1.1 (ii), we see that  $|\mathbf{s}_{ii}| \leq (1 - c_0)/2$  for some small positive constant  $c_0$ . Consequently, since  $S^{(1)}$  is negative definite, we have

$$|\det S^{(1)}| \leq \prod_{i \neq 1} |\mathbf{s}_{ii}| \leq \left(\frac{1 - c_0}{2}\right)^W \quad (10.6)$$

by Hadamard's inequality. Substituting (10.5) and (10.6) into (10.4) yields

$$|\mathcal{I}(\Upsilon_+^b, \Upsilon_-^b, \Upsilon_+^x, \Upsilon_+^x, \Upsilon_S, \mathbb{I}^{W-1})| = O(e^{-cW}) \quad (10.7)$$

for some positive constant  $\delta$ . Hence, we proved the first part of Lemma 5.9. The second part can be proved analogously.  $\square$

In the sequel, we prove Lemma 10.1. We also ignore the factor  $Q(\cdot)$  from the discussion at first.



10.1.  $P(\hat{X}, \hat{B}, V, T)$  in the Type II vicinity. Our aim, is to prove the following lemma.

**Lemma 10.2.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type II vicinity, we have*

$$P(\hat{X}, \hat{B}, V, T) \leq \frac{W^{2+\gamma}\Theta^2}{M} |\det \mathbb{A}_+|^2 \det(S^{(1)})^2. \quad (10.8)$$

*Proof.* We will follow the strategy in Section 9.1. We regard all  $V$ -variables as fixed parameters. Now, we define the function

$$i \equiv i_j(\hat{X}, \hat{B}, T) := |\hat{x}_{j,1}| + |\hat{x}_{j,2}| + |\hat{b}_{j,1}| + |\hat{b}_{j,2}| + |\hat{t}_j|.$$

Then, we recall the representation (9.5) and the definition of  $\Delta_{\ell,j}$  in (9.6). We still adopt the representation (9.9). It is easy to see that in the Type II vicinity, we also have the bound (9.10) for  $\hat{\mathbf{p}}_{\ell,j,\alpha,\beta}$ . The main difference is the first factor of the r.h.s. of (9.5). We expand it around the saddle point as

$$\exp \left\{ -Tr V_j^* \hat{X}_j^{-1} V_j \Omega_j T_j^{-1} \hat{B}_j^{-1} T_j \Xi_j \right\} =: \exp \left\{ -Tr D_+^{-1} \Omega_j D_{\pm}^{-1} \Xi_j \right\} \exp \left\{ -\frac{1}{\sqrt{M}} \hat{\Delta}_j \right\}.$$

We take the formula above as the definition of  $\hat{\Delta}_j$ , which is of the form

$$\hat{\Delta}_j = \sum_{\alpha,\beta=1}^4 \hat{p}_{j,\alpha,\beta} \cdot \omega_{j,\alpha} \xi_{j,\beta},$$

where  $\hat{p}_{j,\alpha,\beta}$  is a function of  $\hat{X}$ ,  $\hat{B}$ ,  $V$  and  $T$ -variables, satisfying

$$\hat{p}_{j,\alpha,\beta} = O(i).$$

Let

$$\hat{\mathbb{H}} = (a_+^{-2} \mathbb{A}_+) \oplus S \oplus (a_-^{-2} \mathbb{A}_+) \oplus S.$$

Recalling the notation in (6.20), we can write

$$-\sum_{j,k} \tilde{\mathbf{s}}_{jk} Tr \Omega_j \Xi_k - \sum_{j=1}^W Tr D_+^{-1} \Omega_j D_{\pm}^{-1} \Xi_j = -\tilde{\Omega} \hat{\mathbb{H}} \tilde{\Xi}'.$$

Now, via replacing  $\Delta_{1,j}$  by  $\hat{\Delta}_{1,j}$ ,  $\hat{\kappa}_j$  by  $i_j$ ,  $\mathbb{H}$  by  $\hat{\mathbb{H}}$  in the proof of Lemma 9.2, we can perform the proof of Lemma 10.2 in the same way. We leave the details to the reader.  $\square$

10.2.  $F(\hat{X}, \hat{B}, V, T)$  in the Type II vicinity. In this section, we will prove the following lemma.

**Lemma 10.3.** *Suppose that the assumptions in Theorem 1.14 hold. In the Type II vicinity, we have*

$$F(\hat{X}, \hat{B}, V, T) = O\left(\frac{\exp\{-(a_+ - a_-)N\eta\}}{(N\eta)^{n+1}}\right). \quad (10.9)$$

*Proof.* Recall the decomposition (9.37). Note that Lemma 9.9 is still applicable. Hence, it suffices to estimate  $\hat{\mathbb{F}}(\hat{X}, V)$ . Now, note that in the Type II vicinity, it is obvious to see that

$$Tr X_j J = O\left(\frac{|\hat{x}_{j,1}| + |\hat{x}_{j,2}|}{\sqrt{M}}\right), \quad \forall j = 1, \dots, W,$$

which implies that

$$\sum_{j=1}^W Tr X_j J = O\left(\frac{\|\hat{\mathbf{x}}_1\|_1 + \|\hat{\mathbf{x}}_2\|_1}{\sqrt{M}}\right) = O\left(\frac{\Theta}{\sqrt{M}}\right).$$

Consequently, we have

$$\exp \left\{ M\eta \sum_{j=1}^W Tr X_j J \right\} = \exp \{ O(\Theta \sqrt{M} \eta) \} = 1 + o(1)$$

by our assumption on  $\eta$ . From (3.19) we can also see that all the other factors of  $f(P_1, V, \hat{X}, X^{[1]})$  are  $O(1)$ . Hence, by the definition (9.36), we have

$$\mathring{\mathbb{F}}(\hat{X}, V) = O(\exp\{-(a_+ - a_-)N\eta\}),$$

which together with Lemma 9.9 yields the conclusion.  $\square$

**10.3. Summing up: Proof of Lemma 10.1.** Analogously, we shall slightly modify the proofs of Lemma 10.2 and Lemma 10.3, in order to take  $Q(\cdot)$  into account. The proof can then be performed in the same manner as Lemma 9.1. We omit the details here.

## 11. PROOF OF THEOREM 1.14

The conclusion for Case 1 is a direct consequence of the discussions in Sections 3.5–10. The proofs of Case 2 and Case 3 can be performed analogously, with slight modifications, which will be stated below.

### • (Case 1)

In this case, by using Lemmas 5.1, 5.6, 5.8 and 5.9, we can get (1.19) immediately.

### • (Case 2)

In this case, we shall slightly modify the discussions in Sections 3.5–10 for Case 1, according to the decomposition of supermatrices in (3.9). Now, at first, in (3.12) and (3.13), for  $A = \check{S}, \check{X}, \check{Y}, \check{\Omega}$  or  $\check{\Xi}$ , we replace  $A_p^{(1)}$  and  $A_q^{(1)}$  by  $A_p^{(1,2)}$  and  $A_q$  respectively, and replace  $A_q^{[1]}$  by  $A_p^{[2]}$ . In addition, in the last three lines of (3.13), we shall also replace  $\tilde{s}_{jq}$  by  $\tilde{s}_{jp}$ , and replace  $\tilde{s}_{pq}$  and  $\tilde{s}_{qp}$  by  $\tilde{s}_{pp}$ , and in the first line, we replace  $\bar{\phi}_{1,q,1}\phi_{1,p,1}\bar{\phi}_{2,p,1}\phi_{2,q,1}$  by  $\bar{\phi}_{1,p,2}\phi_{1,p,1}\bar{\phi}_{2,p,1}\phi_{2,p,2}$ . Then, in (3.14) and (3.15), for  $A = X, Y, \Omega, \Xi, \omega, \xi, \mathbf{w}, y, \tilde{u}, \tilde{v}$  or  $\sigma$ , we replace  $A_q^{[1]}$  by  $A_p^{[2]}$ . With these modifications, it is easy to check the proof in Sections 3.5–10 applies to Case 2 as well. The main point is we can still gain the factor  $1/(N\eta)^{n+1}$  from integral of  $g(\cdot)$  defined in (3.20) (with  $y_q^{[1]}$  and  $\mathbf{w}_q^{[1]}$  replaced by  $y_p^{[2]}$  and  $\mathbf{w}_p^{[2]}$ ). Heuristically, we can go back to (4.6), and replace  $\sigma_q^{[1]}$  by  $\sigma_p^{[2]}$  therein. It is then quite clear the same estimate holds. Consequently, Lemmas 5.1, 5.6, 5.8 and 5.9 still hold under the replacement of the variables described above. Hence, (1.19) holds in Case 2.

### • (Case 3)

Analogously, in this case, we can also mimic the discussions for Case 1 with slight modifications. We also start from (3.12) and (3.13). For  $A = \check{S}, \check{X}, \check{Y}, \check{\Omega}, \check{\Xi}, \omega$  and  $\xi$ , we replace  $A_q^{(1)}$  by  $A_q$ , and replace  $A_q^{[1]}$  by 0. In addition, in the first line of (3.13), we replace  $\bar{\phi}_{1,q,1}\phi_{1,p,1}\bar{\phi}_{2,p,1}\phi_{2,q,1}$  by  $\bar{\phi}_{1,p,1}\phi_{1,p,1}\bar{\phi}_{2,p,1}\phi_{2,p,1}$ . Consequently, after using superbosonization formula, we will get the factor  $(y_p^{[1]}|(\mathbf{w}_p^{[1]}(\mathbf{w}_p^{[1]})^*)_{12}|)^{2n}$  instead of  $(y_p^{[1]}y_q^{[1]}(\mathbf{w}_q^{[1]}(\mathbf{w}_q^{[1]})^*)_{12}(\mathbf{w}_p^{[1]}(\mathbf{w}_p^{[1]})^*)_{21})^n$  in (3.16). Then, for the superdeterminant terms

$$\prod_{k=p,q} \frac{\det(X_k - \Omega_k(Y_k)^{-1}\Xi_k)}{\det Y_k}, \quad \prod_{k=p,q} \frac{y_k^{[1]} \left( y_k^{[1]} - \xi_k^{[1]}(X_k^{[1]})^{-1}\omega_k^{[1]} \right)^2}{\det^2(X_k^{[1]})}.$$

we shall only keep the factors with  $k = p$  and delete those with  $k = q$ . Moreover, we shall also replace  $A_q^{[1]}$  by 0 for  $A = X, Y, \Omega, \Xi, \omega, \xi, \mathbf{w}, y, \tilde{u}, \tilde{v}$  or  $\sigma$  in (3.16). In addition,  $dA^{[1]}$  shall be redefined as the differential of  $A_p^{[1]}$ -variables only, for  $A = X, \mathbf{y}, \mathbf{w}, \mathbf{w}$  and  $\xi$ . One can check step by step that such a modification does not require any essential change of our discussions for Case 1. Especially,

note that our modification has nothing to do with the saddle point analysis on the Gaussian measure  $\exp\{-M(K(\hat{X}, V) + L(\hat{B}, T))\}$ . Moreover, the term  $\mathcal{P}(\cdot)$  in (3.28) can be redefined by deleting the factor with  $k = q$  in the last term therein. Such a modification does not change our analysis of  $\mathcal{P}(\cdot)$ . In addition, the irrelevant term  $\mathcal{Q}(\cdot)$  can also be defined accordingly. Specifically, we shall delete the factor with  $k = q$  in the last term of (3.29) and replace  $A_q^{[1]}$  by 0 for  $A = \Omega, \Xi, \omega, \xi, \mathbf{w}, y$ . It is routine to check that Lemma 6.3 still holds under such a modification. Analogously, we can redefine the functions  $\mathcal{F}(\cdot)$ ,  $f(\cdot)$  and  $g(\cdot)$  in (3.18)-(3.20). Now, the main difference between Case 3 and Case 1 or 2 is that the factor  $(y_p^{[1]} | (\mathbf{w}_p^{[1]} (\mathbf{w}_p^{[1]})^*)_{12} |)^{2n}$  does not produce oscillation in the integral of  $g(\cdot)$  any more. Heuristically, the counterpart of (4.6) in Case 3 reads

$$\begin{aligned} & e^{(a_+ - a_-)N\eta} \int d\mathbf{y}^{[1]} d\mathbf{w}^{[1]} d\nu(Q_1) \cdot g(\hat{B}, T, Q_1, \mathbf{y}^{[1]}, \mathbf{w}^{[1]}) \\ & \sim \int_0^\infty 2t dt \int_{\mathbb{L}} d\sigma_p^{[1]} \cdot e^{-cN\eta t^2 + c_1 e^{-i\sigma_p^{[1]} t}} \sim \frac{1}{N\eta}. \end{aligned}$$

Hence, (1.19) holds for Case 3.

Therefore, we completed the proof of Theorem 1.14.

## 12. FURTHER COMMENTS

In this section, we make some comments on possible further improvements on our results.

- (Comment on how to remove the prefactor  $N^{C_0}$  in (1.19))

As mentioned in Remark 1.15, we have used  $N^{C_0}$  to replace  $M\Theta^2 W^{C_0}/(N\eta)^\ell$ . However, the latter is also artificial. It can be improved to some  $n$ -dependent constant  $C_n$  via a more delicate analysis on  $A(\cdot)$ , i.e. the integral of  $\mathcal{P}(\cdot)\mathcal{Q}(\cdot)\mathcal{F}(\cdot)$ . Such an improvement stems from the cancellation in the Gaussian integral. At first, a finer analysis will show that the factor  $\mathcal{Q}(\cdot)$  can really be ignored, in the sense that it does not play any role in the estimate of the order of  $\mathbb{E}|G_{ij}(z)|^{2n}$ . Hence, for simplicity, we just focus on the product  $P(\cdot)F(\cdot)$  instead of  $A(\cdot)$ . Then, we go back to Lemma 9.2 and Lemma 9.7. Recall the decomposition (9.37). A more careful analysis on  $F(\cdot)$  leads us to the following expansion, up to the subleading order terms of the factors  $\mathring{\mathbb{G}}(\cdot)$  and  $\mathring{\mathbb{F}}(\cdot)$ ,

$$\begin{aligned} F(\cdot) = \mathring{\mathbb{G}}(\cdot)\mathring{\mathbb{F}}(\cdot) & \sim \frac{1}{(N\eta)^{n+2}} \left( 1 + \frac{M\eta}{\sqrt{M}} \sum_{j=1}^W l_j(\hat{x}_{j,1}, \hat{x}_{j,2}, \hat{v}_j) + \cdots \right) \\ & \times \left( 1 + \frac{M\eta}{\sqrt{M}} \sum_{j=1}^W l'_j(\hat{b}_{j,1}, \hat{b}_{j,2}, \hat{t}_j) + \cdots \right), \end{aligned} \quad (12.1)$$

where  $l_j(\cdot)$ 's and  $l'_j(\cdot)$ 's are some linear combinations of the arguments. Analogously, we shall write down the leading order term of  $P(\cdot)$  in terms of  $\hat{\mathbf{x}}, \hat{\mathbf{b}}, \hat{\mathbf{t}}$  and  $\hat{\mathbf{v}}$  explicitly, instead of bounding it crudely by using (9.24). Then it can be seen that the leading order term of  $P(\cdot)$  is a linear combination of  $\hat{x}_{j,1}\hat{x}_{k,2}$ ,  $\hat{b}_{j,1}\hat{b}_{k,2}$ ,  $\hat{x}_{j,\alpha}\hat{b}_{k,\beta}$ ,  $v_{j,\alpha}\tau_{k,\beta}$  for  $j, k = 1, \dots, W$  and  $\alpha, \beta = 1, 2$ , in which all the coefficients are of order  $1/M$ . Observe that the Gaussian integral in (8.47) will kill the linear terms. Consequently, in the expansion (12.1), the first term that survives after the Gaussian integral is actually

$$\frac{1}{(N\eta)^{n+2}} \cdot \frac{M\eta}{\sqrt{M}} \sum_{j=1}^W l_j(\hat{x}_{j,1}, \hat{x}_{j,2}, \hat{v}_j) \cdot \frac{M\eta}{\sqrt{M}} \sum_{j=1}^W l'_j(\hat{b}_{j,1}, \hat{b}_{j,2}, \hat{t}_j). \quad (12.2)$$

Replacing  $A(\cdot)$  by the product of the leading order term of  $P(\cdot)$  and (12.2) in the integral (8.47) and taking the Gaussian integral over  $\mathbf{c}$ ,  $\mathbf{d}$ ,  $\boldsymbol{\tau}$  and  $\mathbf{v}$ -variables yield the true order  $1/(N\eta)^n$ , without additional  $N$ -dependent prefactors.

- (Comment on the restriction  $|E| \leq \sqrt{2} - \kappa$ )

This restriction is used in several places, we mention the two most important ones.

The most critical issue is the term  $\ell_S(\hat{B}, T)$  defined in (5.6). A direct consequence of  $|E| \leq \sqrt{2} - \kappa$  is that  $\operatorname{Re}(b_{j,1} + b_{j,2})(b_{k,1} + b_{k,2}) \geq 0$  for  $\mathbf{b}_1 \in \Gamma^W$  and  $\mathbf{b}_2 \in \bar{\Gamma}^W$ , thus  $\operatorname{Re} \ell_S(\hat{B}, T) \geq 0$ . However, once  $|E| > \sqrt{2}$ ,  $\operatorname{Re} \ell_S(\hat{B}, T)$  can be negative for  $\mathbf{b}_1 \in \Gamma^W$  and  $\mathbf{b}_2 \in \bar{\Gamma}^W$ . Consequently, the measure  $\exp\{-M\ell_S(\hat{B}, T)\}$  is not well defined, considering the domain of the  $\mathbf{t}$ -variables is not compact. Actually, such a problem is unavoidable if we independently deform the contours of  $b_{j,1}$ 's (resp.  $b_{j,2}$ 's) from  $\mathbb{R}_+$  to any contour passing through the saddle point  $a_+$  (resp.  $a_-$ ), starting from the 0, since one can always choose  $b_{j,1}$  and  $b_{k,1}$  to be  $a_+$  while  $b_{j,2}$  and  $b_{k,2}$  to be very close to 0, such that  $\operatorname{Re}(b_{j,1} + b_{j,2})(b_{k,1} + b_{k,2}) < 0$  in case  $|E| > \sqrt{2}$ . A possible way to solve this problem is to change the variables  $(b_{j,1}, b_{j,2})$  to  $(r_{j,1}, r_{j,2})$ , defined by  $r_{j,1} = (b_{j,1} + b_{j,2})/2$  and  $r_{j,2} = (b_{j,1} - b_{j,2})/2$ , and perform the saddle point analysis with respect to the latter. It is not difficult to calculate that the saddle point of  $L(\hat{B}, T)$  restricted on these  $\mathbf{r}$ -variables is  $r_{j,1} = \sqrt{4 - E^2}$  and  $r_{j,2} = \mathbf{i}E$ . To guarantee the positivity of  $\ell_S(\hat{B}, T)$ , it suffices to choose a contour for  $r_{j,1}$  passing through  $\sqrt{4 - E^2}$ , staying in the sector  $\tilde{\mathbb{K}} := \{\omega \in \mathbb{C} : |\arg \omega| < \pi/4\}$ . However, the explicit way to choose the contours for  $r_{j,1}$  and  $r_{j,2}$  and the subsequent analysis in terms of  $\mathbf{r}$ -variables instead of  $\mathbf{b}$ -variables would be more involved and we leave it to future work.

Another point where  $|E| < \sqrt{2} - \kappa$  is used is the proof of Lemma 5.9 in Section 10. In (10.3), we used  $\operatorname{Re} \mathbb{A}_+^v \geq I$ , which is a consequence of  $\operatorname{Re} a_+^2 = (2 - E^2)/2 \geq 0$  and (5.23), see the definition of  $\mathbb{A}_+^v$  in (8.35). If  $|E| > \sqrt{2}$ , the bound in (10.3) should be weakened to  $\operatorname{Re}(\tilde{\mathbf{x}}' \mathbb{A}_+^v \tilde{\mathbf{x}}) \geq c(E) \|\tilde{\mathbf{x}}\|_2^2$  for some  $E$ -dependent constant  $c(E) < 1$ , which is uniform on  $V_j \in \mathring{U}(2)$  for all  $j = 2, \dots, W$ . This modification, however, is not critical. An easy way to remedy this situation is to impose an additional condition  $\det S^{(1)} \leq (c(E) - c)^W$  for arbitrarily small constant  $c$ , to guarantee the estimate (10.7).

## REFERENCES

- [1] Aizenman, M., and Molchanov, S.: Localization at large disorder and at extreme energies: an elementary derivation, *Commun. Math. Phys.* **157**, 245–278 (1993)
- [2] Anderson, P.: Absences of diffusion in certain random lattices, *Phys. Rev.* **109**, 1492–1505 (1958)
- [3] Bunder J.E., Efetov K.B., Kravtsov V.E., Yevtushenko O.M. and Zirnbauer M.R.: Superbosonization formula and its application to random matrix theory. *J. Stat. Phys.* **129** (5–6), 809–832 (2007)
- [4] Disertori, M., Pinson, H., Spencer, T.: Density of states for random band matrices. *Commun. Math. Phys.* **232**(1), 83–124 (2002)
- [5] Disertori, M., Spencer, T.: Anderson localization for a supersymmetric sigma model. *Commun. Math. Phys.* **300**(3), 659–671 (2010)
- [6] Disertori, M., Spencer, T., Zirnbauer, M.R.: Quasi-diffusion in a 3D supersymmetric hyperbolic sigma model. *Commun. Math. Phys.* **300**(2), 435–486 (2010)
- [7] Efetov, K.: Supersymmetry in disorder and chaos. Cambridge University Press, 1997.
- [8] Ellis, R. B.: Discrete Green's functions for products of regular graphs. arXiv preprint math/0309080, (2003)
- [9] Erdős, L., Knowles, A., Yau, H.-T.: Averaging fluctuations in resolvents of random band matrices. *Ann. Henri Poincaré* **14**(8), 1837–1926 (2013)
- [10] Erdős, L., Knowles, A., Yau, H.-T., and Yin, J.: Spectral statistics of Erdős-Rényi graphs I: local semicircle law. *Ann. Probab.* **41**(3B), 2279–2375 (2013)
- [11] Erdős, L., Knowles, A., Yau, H.-T., Yin, J.: Delocalization and Diffusion Profile for Random Band Matrices. *Commun. Math. Phys.* **323**, 367–416 (2013)

- [12] Erdős, L., Knowles, A., Yau, H.-T., Yin, J.: The local semicircle law for a general class of random matrices. *Electron. J. Probab.* **18**(59), 1–58 (2013)
- [13] Erdős, L., Schlein, B., Yau, H.-T.: Local semicircle law and complete delocalization for Wigner random matrices. *Commun. Math. Phys.* **287**, 641–655 (2009)
- [14] Erdős, L., Yau, H.-T., Yin, J.: Bulk universality for generalized Wigner matrices. *Probab. Theory and Related Fields* **154**(1-2), 341–407 (2012)
- [15] Fröhlich, J., Spencer, T.: Absence of diffusion in the Anderson tight binding model for large disorder or low energy, *Commun. Math. Phys.* **88**, 151–184 (1983)
- [16] Fyodorov, Y.V. and Mirlin, A.D.: Scaling properties of localization in random band matrices: a  $\sigma$ -model approach. *Phys. Rev. Lett.* **67**, 2405–2409 (1991)
- [17] Littelmann, P., Sommers, H.J., Zirnbauer, M.R.: Superbosonization of invariant random matrix ensembles. *Commun. Math. Phys.* **283**(2), 343–395 (2008)
- [18] Shcherbina, T.: On the second mixed moment of the characteristic polynomials of 1D band matrices. *Commun. Math. Phys.* **328**(1), 45–82 (2014)
- [19] Shcherbina, T.: Universality of the local regime for the block band matrices with a finite number of blocks. *J. Stat. Phys.* **155**(3), 466–499 (2014)
- [20] Shcherbina, T.: Universality of the second mixed moment of the characteristic polynomials of the 1D band matrices: real symmetric case. arXiv:1410.3084, 2014.
- [21] Schäfer, L., Wegner, F.: Disordered system with  $n$  orbitals per site: Lagrange formulation, hyperbolic symmetry, and Goldstone modes. *Z. Phys. B* **38**, 113–126 (1980)
- [22] Schenker, J.: Eigenvector localization for random band matrices with power law band width. *Commun. Math. Phys.* **290**, 1065–1097 (2009)
- [23] Sodin, S.: An estimate for the average spectral measure of random band matrices, *J. Stat. Phys.* **144** (2011), no. 1, 46–59.
- [24] Sodin, S.: The spectral edge of some random band matrices. *Ann. of Math.* **172**(3), 2223–2251 (2010)
- [25] Spencer, T.: SUSY statistical mechanics and random band matrices. Quantum many body systems. Springer Berlin Heidelberg, 125–177 (2012)
- [26] Spencer, T., Zirnbauer, M.R.: Spontaneous symmetry breaking of a hyperbolic sigma model in three dimensions. *Commun. Math. Phys.* **252**(1-3), 167–187 (2004)
- [27] Tao, T., Vu, V.: Random matrices: universality of local eigenvalue statistics. *Acta Math.* **206**(1), 127–204 (2011)
- [28] Wegner, F.J.: Disordered system with  $n$  orbitals per site:  $n \rightarrow \infty$  limit, *Phys. Rev. B* **19**, 783–792 (1979)
- [29] Wigner, E.: Characteristic vectors of bordered matrices with infinite dimensions. *Ann. of Math.* **62**, 548–564 (1955)

IST AUSTRIA, AM CAMPUS 1, 3400 KLOSTERNEUBERG, AUSTRIA

E-mail address: maomie2007@gmail.com

IST AUSTRIA, AM CAMPUS 1, 3400 KLOSTERNEUBERG, AUSTRIA

E-mail address: laszlo.erdos@ist.ac.at